## Leveraging Large Language Models for

#### **Automated Chart Summarization**

#### Romain Bazin



# Department of Electrical & Computer Engineering McGill University, Montreal

October 1st, 2024

A thesis submitted to McGill University in partial fulfillment of the requirements of the degree of

Master of Science

© 2024 Romain Bazin

#### Abstract

Access to graphical information on the internet remains a significant barrier for blind or visually impaired individuals, particularly when it comes to data visualizations like charts. This thesis explores how recent advancements in artificial intelligence can enhance accessibility through innovations in chart summarization, the process of automatically extracting information from a chart and compiling it into a textual summary intelligible to screen readers.

The thesis focuses on two main areas of research. First, it applies the emerging architecture of large language model agents to the task of chart summarization, a novel application in this domain. This approach combines recent advancements in chart information extraction with the reasoning and planning capabilities of large language models. By leveraging natural language processing technologies, it reduces the need for curation of annotated datasets traditionally required to train vision AI models. An implementation of this agent-based approach is developed and evaluated, demonstrating its effectiveness in generating chart summaries.

Abstract

Second, the thesis addresses the challenge of evaluating chart summarization systems. Adapting approaches from the field of text summarization, it introduces a framework of criteria to assess various aspects of chart summaries quality. This framework serves to identify user preferences while enabling accurate measurement and classification of various chart summarization systems. The framework is assessed through a comparative study with human evaluators, providing insights into its effectiveness for evaluating automated chart summarization systems.

The results of this research will hopefully inform future developments in chart summarization techniques and user-adaptive accessibility solutions.

## Abrégé

L'accès aux informations visuelles sur internet demeure un obstacle majeur pour les personnes aveugles ou malvoyantes, en particulier concernant les visualisations de données telles que les graphiques. Cette thèse explore comment les récents progrès en intelligence artificielle peuvent améliorer l'accessibilité grâce à des innovations dans la synthèse de graphiques, un processus d'extraction automatique d'informations à partir d'un graphique et de compilation en un résumé textuel intelligible pour les lecteurs d'écran.

La thèse se concentre sur deux principaux axes de recherche. Premièrement, elle applique l'architecture émergente des agents basés sur les grands modèles de langage à la tâche de synthèse de graphiques, une application novatrice dans ce domaine. Cette approche combine les avancées récentes en extraction d'informations graphiques avec les capacités de raisonnement et de planification des grands modèles de langage. En s'appuyant sur les technologies de traitement du langage naturel, elle réduit le besoin de curation des jeux de données annotés traditionnellement nécessaires pour entraîner les modèles de vision par intelligence artificielle. Une implémentation de cette approche basée

Abrégé iv

sur des agents est développée et évaluée, démontrant son efficacité dans la génération de résumés de graphiques.

Deuxièmement, la thèse aborde le défi de l'évaluation des systèmes de synthèse de graphiques. En adaptant des approches du domaine de la synthèse de texte, elle introduit un des critères pour évaluer divers aspects de la qualité des résumés de graphiques. Ces critères permettent d'identifier les préférences des utilisateurs tout en permettant une mesure et une classification précises des différents systèmes de synthèse de graphiques. Cette méthode est appliquée dans une étude comparative avec des évaluateurs humains, révélant son utilité pour comparer les modèles et identifier les préférences entre différents groupes d'utilisateurs.

Les résultats de cette recherche pourront, nous l'espérons, éclairer les futurs développements dans les techniques de synthèse de graphiques et les solutions d'accessibilité adaptatives pour les utilisateurs.

## Acknowledgements

I extend my deepest gratitude to my supervisor, Professor Jeremy Cooperstock, for his invaluable mentorship, support, and provision of essential infrastructure throughout this journey. His encouragement to follow my research interests has been fundamental to my work.

My heartfelt thanks go to my colleagues at the Shared Reality Lab for their insightful suggestions and advice. Special appreciation is due to Cyan Kuo, whose consistent guidance has been instrumental to my progress.

I am profoundly grateful to my family for their unwavering support. My parents fostered a nurturing environment that paved the way for my achievements, while my sister's affectionate nature has consistently provided comfort and strength. I extend my heartfelt thanks to Cyril, my best friend, whose faithful friendship has been a constant source of encouragement throughout the years. Finally, I reserve my deepest appreciation for Jade, whose enduring love and support have been essential to my journey and success.

This journey would not have been possible without the collective support and

encouragement of each one mentioned, and for that, I am eternally grateful.

I would also like to acknowledge the helpful assistance of large language models in refining this thesis. In particular, Claude 3.5 Sonnet has been valuable in helping me rephrase and improve my writing, contributing to the clarity and coherence of this work. While the ideas and research are entirely my own, this tool has enhanced the presentation of my thoughts and findings.

## Contents

1	Intr	Introduction		
2	Bac	ackground		5
	2.1	Advar	nces in Computer Vision and Natural Language Processing for Chart	
		Under	estanding	5
		2.1.1	The Importance of Chart Understanding in Improving Accessibility for	
			Blind or Visually Impaired People	5
		2.1.2	Evolution of Image Understanding and Chart Understanding Methods	8
		2.1.3	Dataset Scarcity in Chart Summarization	10
		2.1.4	Opportunities from Chart Visual Question Answering	11
		2.1.5	Conclusion	13
	2.2	Large	Language Model Agents and Machine Learning Models as Tools	14
		2.2.1	Large Language Models as AI Agents	14
			2.2.1.1 Brain Module of LLM-based AI Agents	1/

Contents

			2.2.1.2 Perception Module of LLM-based AI Agents	15
			2.2.1.3 Action Module of LLM-based AI Agents	17
		2.2.2	Conclusion	21
	2.3	Limits	s of Chart Summarization Evaluation Methods	22
		2.3.1	Current Evaluation Methods in Chart Summarization and Their	
			Limitations	22
		2.3.2	Inspiration and Opportunities from Text Summarization Evaluation $$ .	24
		2.3.3	Conclusion	26
3	Des	ign of	Large Language Model Agent for Chart Summarization	27
	3.1	Archit	tecture Overview	28
	3.2	Modu	les Implementation	32
		3.2.1	Perception Module	32
		3.2.2	Brain Module	33
			3.2.2.1 Prompt Engineering Techniques	33
			3.2.2.2 Tool Learning and Utilization	34
			3.2.2.3 Memory Mechanism	36
			3.2.2.4 Summary Generation	37
		3.2.3	Action Module	40
			3.2.3.1 Tool Interaction Protocol	40
			3.2.3.2 Tools List	42

ix

	3.3	Discussion	44
		3.3.1 Design Choices and Trade-offs	44
		3.3.2 Future Directions and Implications	46
4	Sys1	tem Evaluation	47
	4.1	Dataset	49
	4.2	Baselines	52
	4.3	Quality Criteria Evaluation Framework	54
	4.4	User Study	57
		4.4.1 Participant Selection and Categorization	58
		4.4.2 Materials Preparation	60
		4.4.3 Data Collection Procedures	63
	4.5	Results	65
		4.5.1 User Preferences	66
		4.5.2 Quality Criteria	71
	4.6	Discussion	78
5	Disc	cussion	82
	5.1	Summary of Themes and Key Findings	82
	5.2	Implications and Practical Applications	84
	5.3	Limitations	85

	5.4	Future Research Directions	85
	5.5	Conclusion	86
$\mathbf{A}$	Exp	ertise Questionnaire	101
В	Sum	amary Generation Example	106

## List of Figures

3.1	Summarization System Architecture Overview	29
4.1	Examples of Charts from Pew Chart2Text Subset	49
4.2	Topic and Chart Type Distributions in Pew Chart2Text Subset	50
4.3	Comparison of original and modified charts	60
4.4	Ranking Distribution of Chart Summarization Models	67
4.5	Model Performance Comparison Across Quality Criteria	72
B 1	Percentage of U.S. adults' cash purchases by income level	106

## List of Tables

4.1	Nemenyi Post-hoc Test Results for pair-wise Model Comparisons	68
4.2	ART ANOVA Results for Chart Summarization Study Factors	69
4.3	Expert vs. Novice Model Rankings Comparison	70
4.4	Friedman Test Results for Quality Criteria Across Models	73
4.5	Nemenyi Test Results for Model Comparisons Across Quality Criteria	74
4.6	Criteria-Ranking Correlation Analysis	77

## List of Acronyms

**AI** Artificial Intelligence.

**BVIP** Blind or Visually Impaired People.

**COT** Chain of Thought.

CVQA Chart Visual Question Answering.

 $\mathbf{GPU} \qquad \text{Graphics Processing Unit.}$ 

ICL In-Context Learning.

 $\mathbf{LLM} \qquad \text{Large Language Model}.$ 

**QC** Quality Criteria.

**TOT** Tree of Thoughts.

**VFM** Visual Foundation Model.

**VIT** Vision Transformer.

VLLM Vision Large Language Model.

#### Contributions

This thesis applies LLM-based agent architecture to chart summarization, addressing dataset scarcity in this domain. The approach combines chart information extraction techniques with large language model reasoning capabilities, offering a modular solution for generating summaries without relying on annotated datasets. While not outperforming the closed-source GPT-4-Vision, our method demonstrates competitive results, particularly in linguistic aspects, and surpasses the previous open-source state-of-the-art, Unichart. This contribution provides academia and open-source developers a pathway to potentially match proprietary model performance through improved information extraction and advanced open-source LLMs like LLAMA-3.

Our research also identifies differences in chart summary preferences between expert and novice users. Experts value a balance of linguistic and data-centric aspects, while novices prioritize relevance. These findings contribute to understanding diverse user needs in chart summarization and indicate potential for developing adaptive systems tailored to varying levels of expertise.

## Chapter 1

#### Introduction

Artificial intelligence (AI) is significantly impacting industries and creating new possibilities across diverse domains. From healthcare and education to transportation and beyond, AI is driving substantial innovation and efficiency improvements. Two subfields at the forefront of this development are computer vision and natural language processing. Computer vision focuses on enabling machines to interpret and understand visual information, while natural language processing aims to teach computers to comprehend and communicate in human language.

The advancements in computer vision and natural language processing have created new opportunities for AI-driven accessibility solutions. An important challenge at the intersection of these fields is the automatic summarization of visual data representations, particularly charts and graphs. While graph summarization has seen significant progress, largely due

1. Introduction 2

to advancements in image captioning techniques and the availability of large-scale datasets such as Common Objects in Context (COCO) [1] and Flickr30k [2], the summarization of charts presents unique challenges that warrant further research .

Charts are a prevalent means of conveying complex data and relationships in a visual format, used extensively in scientific publications, business reports, and online media. Unlike graphs, which primarily represent relationships between entities, charts often combine visual elements, textual information, and numerical data values. This multimodal nature makes chart summarization challenging, requiring advanced AI techniques to interpret and synthesize information accurately. This challenge is notable in the context of accessibility for blind and visually impaired (BVI) individuals. While the text surrounding a chart may provide some context, it often doesn't fully capture the information presented in the visual representation. Despite recent progress in chart summarization techniques, many current solutions still struggle to fully convey the depth of information in a format accessible to BVI users [3, 4]. This gap in accessibility leaves BVI users potentially excluded from accessing important information, limiting their ability to fully engage with data-rich content.

Generating comprehensive and accurate summaries of charts is a challenging AI task, requiring an understanding of diverse chart types, data encodings, and represented topics. A key bottleneck in developing effective chart summarization systems is the need for large, high-quality training datasets. The intricacy and variety of charts demand substantial amounts of

1. Introduction 3

labeled examples to train models that can generalize well to the vast space of unseen charts. However, collecting such datasets through manual human annotation is a time-consuming and expensive process. Despite recent progress in chart summarization techniques, many current solutions still struggle to fully convey the depth of information in a format accessible to BVI users [3, 4]. This presents a significant challenge for progress in this area, as the cost and effort required to build suitable training data are often considerable.

The first key objective of this thesis is to develop a novel paradigm for chart summarization that harnesses the power of large language models (LLMs). By leveraging the extensive general knowledge and reasoning capabilities of LLMs, this approach aims to reduce the reliance on large labeled training datasets, thus overcoming a critical limitation of prior work.

While prior research has made strides in extracting information from charts, less attention has been devoted to the equally important task of synthesizing this information into coherent, insightful summaries that align with user preferences. The lack of well-defined, measurable attributes for assessing summary quality has hindered progress in understanding what makes a chart summary truly effective and valuable to users. To address this gap, the second key objective of this thesis is to establish a framework of quantitative criteria for evaluating the quality of chart summaries. Inspired by recent advancements in text summarization evaluation, we propose a set of criteria, each rated on a 1-5 scale, that capture essential aspects of summary quality.

1. Introduction 4

To address these objectives and present our research, this thesis is organized into five chapters

Chapter 2: Background reviews the current state of chart understanding. It examines recent developments in computer vision and natural language processing techniques, while also discussing the limitations of existing chart summarization approaches, with particular attention to the scarcity of high-quality datasets.

Chapter 3: Design of LLM Agent for Chart Summarization introduces the primary contribution of this thesis: an LLM-based agent architecture for chart summarization. It provides a detailed description of the brain, perception, and action modules that constitute this system.

Chapter 4: System Evaluation assesses the performance of the proposed approach against relevant baselines. It introduces a quality criteria framework adapted from text summarization and presents a user study demonstrating its efficacy in capturing nuanced differences in summary quality across models and user groups.

Chapter 5: Discussion analyzes the key findings on LLM-based chart summarization. It considers the limitations of the proposed approach and identifies potential future research directions. The chapter concludes by revisiting the initial objectives and exploring their implications for improving access to visual data.

## Chapter 2

## Background

- 2.1 Advances in Computer Vision and Natural

  Language Processing for Chart Understanding
- 2.1.1 The Importance of Chart Understanding in Improving

  Accessibility for Blind or Visually Impaired People

Chart understanding is a subfield of image understanding focused on automatically extracting and interpreting information from charts to facilitate user comprehension. It has a long history of research interest due to its importance in making complex data more understandable. Early work included knowledge-based systems like WERP in the 1980s [5], which generated weather reports from charts, and efforts in the 2000s [6] to associate text

and graphics in scientific charts for more comprehensive semantic understanding.

Chart understanding is important for increasing accessibility for blind or visually impaired people (BVIP), who face significant challenges when accessing visual information in charts. While BVIP rely on assistive technologies like screen readers to navigate digital content, these tools often struggle to convey the full depth of insights and nuances contained within complex visual representations. The surrounding text may provide some context, but it rarely captures all the information encoded in the visual elements, such as trends, patterns, data points relationships, and intended message. Consequently, BVIP are frequently excluded from accessing critical information, hindering their ability to fully engage with the content and participate in data-driven discussions [7].

To address this accessibility gap, researchers have explored various approaches to extracting meaningful information from charts and presenting it in alternative formats that are more accessible to BVIP. These efforts include developing techniques for automatically parsing and analyzing chart images to identify key components such as axes, labels, and data points; a task commonly known as chart understanding [8, 9, 10, 11, 3]. By extracting this structured information, it becomes possible to generate alternative renderings of charts that convey the same insights through non-visual modalities, such as audio signals, vibrotactile feedback or force feedback haptics [12, 13]. The ultimate goal is to create a seamless and engaging experience for BVIP by combining these alternative renderings into a coherent and informative format [14, 15, 16].

Image-to-text chart understanding is a promising approach to bridging the accessibility gap for BVIP by translating chart images into textual representations. It encompasses three main tasks: chart visual question answering (CVQA), chart derendering, and chart summarization. CVQA involves providing accurate answers to natural language questions about a chart, requiring an understanding of both visual elements and textual content. Chart derendering, also known as chart-to-table, converts a chart image into a structured data table by identifying and extracting components such as bars, lines, legends, and axes. Chart summarization, the primary focus of this thesis, generates concise textual summaries that capture the key insights and salient information in a chart by analyzing visual and textual elements to produce a narrative highlighting the most important aspects of the data.

Among these tasks, chart summarization has received less attention in the research community but provides significant value for BVIP. By generating informative, well-articulated text summaries encapsulating the main takeaways of a chart, chart summarization enables BVIP to access and comprehend the essential information conveyed in visual representations. This thesis aims to advance chart summarization methods by developing new methods in summary generation and evaluation.

The methods employed in image-to-text chart understanding, particularly in chart summarization, draw inspiration from advancements in the field of image understanding [17]. The next section will present the evolution of image captioning and chart summarization methods, exploring the architectures and techniques that have been

adapted from image understanding to tackle the challenges of extracting and conveying information from charts in a textual format.

# 2.1.2 Evolution of Image Understanding and Chart Understanding Methods

A comprehensive survey on image captioning methods, a task closely related to chart summarization, provides theoretical insights into the architectures used in image understanding systems [18]. The image captioning pipeline typically consists of two main components: a visual detector for extracting factual information from visual elements, and a description generator for producing descriptions, answers, or tables based on the extracted visual information.

Over the years, the implementation of these components has evolved significantly. Early machine learning-based approaches included template-based methods [19], which relied on predefined templates and statistical language models to generate captions using object detectors. This approach was applied to chart summarization, using a combination of neural networks for visual detection and three-part templates for description generation [10], consisting of a premier (describing basic shape and factual information), a core (conveying the intended meaning of the chart), and a wrapper (providing additional details for certain chart types). However, these methods were limited by their reliance on manually defined features and templates, restricting their generalization capabilities.

Retrieval-based methods focused on ranking existing captions or images based on their similarity to the input, selecting the most relevant ones from a predefined dataset [20]. The limitations here were the dependency on dataset diversity and size for caption relevance and novelty, and the reliance on large labeled datasets with simple, short captions, making them unsuitable for chart summarization.

In recent years, there has been a paradigm shift towards deep learning-based approaches for image captioning. Encoder-decoder models, inspired by machine translation systems, encode visual information using convolutional neural networks (CNNs) and decode it into textual descriptions with recurrent neural networks (RNNs) [21, 22, 23]. This end-to-end approach allows for more flexible and contextually relevant captions by learning features directly from data. It has been applied to chart summarization in works like STLCQA [24]. Attention-based models further refine the encoder-decoder architecture by enabling the model to focus on specific parts of the image while generating each word of the caption, mimicking human visual attention [25]. This approach leads to more detailed and accurate descriptions and has been extensively used for CVQA and chart derendering models, such as ChartQA [26], MatCha [11], and DePlot [27] models, but also in chart summarization as in the Chart-to-Text [4] and UniChart models [3].

#### 2.1.3 Dataset Scarcity in Chart Summarization

While Deep Learning methods have become the dominant approach in image captioning, CVQA, and chart derendering, they require large amounts of high-quality labeled data to reach their full potential. Unfortunately, as highlighted in a review on chart classification and captioning [17], the main limitation in applying these techniques to chart summarization is the scarcity of such datasets. Fredriksson et al. (2022) [28] found that annotation costs vary significantly with task complexity and required domain expertise. More constrained tasks, such as verification or simple identification, require less time and resources than open-ended, creative tasks. This particularly impacts chart summarization, which requires creating coherent, multi-sentence descriptions, compared to CVQA's brief, focused answers. With annotation time increasing non-linearly with task length, and validation costs also rising with complexity, creating high-quality summarization datasets becomes prohibitively expensive compared to question-answering ones. Recent efforts like Chart-to-Text [4] have aimed to create a large labeled dataset with over 35k chart summaries. However, the dataset suffers from issues with the label summaries containing information not present in the charts, leading to models trained on it exhibiting severe hallucination problems [29]. These models generate coherent but false information about new charts, undermining their real-world usability.

To address dataset scarcity, the creators of the UniChart [3] model proposed crafting a synthetic dataset from a harmonized corpus of real-world charts. They collected 627,000

charts from various sources, including online databases, publications, and existing datasets. For charts with available data tables, they utilized existing datasets or extracted data from SVG elements. They also performed data augmentation by creating new charts from public data tables using visualization libraries. For charts without data tables, they employed object detection and optical character recognition (OCR) to extract textual and visual elements. To generate summaries for the charts, they used a knowledge distillation approach. First, they used a large language model, InstructGPT [30], to generate summaries for a subset of charts based on extracted textual information, which were then reviewed and validated by human annotators. These summaries were used to fine-tune another model, Flan-T5 XL [31], which subsequently generated summaries for the larger dataset. The model trained on the larger dataset, using an attention-based encoder-decoder architecture, achieved state-of-the-art results in chart summarization. However, this large-scale dataset was not made public, so the problem of dataset scarcity persists.

#### 2.1.4 Opportunities from Chart Visual Question Answering

In contrast to the scarcity of chart summarization datasets, the fields of CVQA and chart derendering have thrived due to the emergence of multiple high-quality datasets that are easier to create. Notable CVQA datasets include DVQA [32], FigureQA [33], LeafQA [34], PlotQA [35], ChartQA [26], and OpenCQA [36]. The availability of these datasets has led to a proliferation of research on CVQA and chart derendering, ranging from early systems

combining manual and automated chart specification for exploration and question answering [37], to more recent work like MatCha [11] and DePlot [27].

MatCha, enhances visual language pretraining with math reasoning and chart derendering tasks. MatCha is built upon an image-to-text transformer architecture and learns to extract information from charts by predicting the underlying data table or the code used to render the chart. Additionally, it is pretrained on math reasoning tasks using datasets like MATH and DROP, enabling it to perform numerical reasoning on the extracted information. This pretraining allows MatCha to achieve state-of-the-art results on several CVQA benchmarks.

DePlot extends the chart derendering capability of MatCha by fine-tuning the model on a larger dataset of chart-table pairs, achieving even greater accuracy in extracting structured data from charts. Furthermore, DePlot improves upon MatCha visual language reasoning capabilities by proposing a modular approach that separates the task into two steps: (1) chart-to-table conversion using the DePlot model, and (2) reasoning over the extracted table using a LLM. In the first step, the DePlot model converts the chart image into a structured table format. The output of this model is then fed into an LLM, which is prompted to answer questions or perform reasoning tasks based on the extracted table. By leveraging the few-shot learning capabilities of LLMs, DePlot achieves considerable results on CVQA tasks with minimal human-labeled examples, highlighting the promise of this modular approach in scenarios where labeled data is scarce [27].

Following these works, the UniChart model utilized MatCha pretraining scheme and DePlot's fine-tuning method on their private large-scale dataset to reach even higher accuracy on CVQA and chart derendering tasks.

This presents an opportunity to leverage the advances and resources in CVQA to improve chart summarization. By integrating highly accurate CVQA and derendering models into an automated summarization workflow, we can potentially craft more accurate, coherent, and comprehensive summaries, even in the face of limited labeled summarization datasets.

#### 2.1.5 Conclusion

In conclusion, while chart understanding has progressed significantly, chart summarization remains challenging due to the scarcity of high-quality labeled datasets. Models from related fields like CVQA and chart derendering offer promising building blocks, but require sophisticated integration to generate coherent summaries. Recent advances in LLMs present an opportunity to bridge this gap through their strong language and cognitive capabilities. The following section introduces an LLM-based agent architecture designed to use these components for robust chart summarization.

## 2.2 Large Language Model Agents and Machine Learning Models as Tools

#### 2.2.1 Large Language Models as AI Agents

AI agents are artificial entities that perceive their environment, make decisions, and take actions to achieve specified goals, with recent advancements in LLMs enabling increasingly capable implementations [38]. These LLM-based agents rely on a modular architecture, consisting of a central brain module that interacts with perception and action modules, to understand and interact with their environment effectively, making them particularly promising for complex tasks like chart summarization.

#### 2.2.1.1 Brain Module of LLM-based AI Agents

The brain module, centered around an LLM, serves as the hub for various cognitive capabilities, including natural language interaction, knowledge representation, memory management, reasoning, planning, and generalization [38]. LLMs excel at engaging in multi-turn conversations, generating coherent text, and comprehending the intentions behind language to a certain extent [39]. Instruction tuning enables LLM-based agents to generalize to new tasks without additional training, a capability known as zero-shot [40], or few-shot generalization [41]. This flexibility is particularly valuable for tasks like chart summarization, which may not have been included in the pretraining or fine-tuning data.

To enhance the reasoning and planning capabilities of LLM-based agents, techniques such as chain-of-thought [42] prompting, self-consistency [43] methods, and task decomposition have been developed. Chain-of-thought prompting encourages the model to explicitly establish a thought process before generating answers, while self-consistency methods enable the model to explore multiple reasoning paths for optimal decision making. Task decomposition breaks down complex problems into more manageable sub-tasks, making problem-solving more efficient and effective.

However, challenges remain in understanding implicit information [44] and mitigating hallucinations [45], where the model generates inconsistent or factually incorrect content. Memory management is also a crucial aspect of the brain module, as the amount of memory an LLM can process at once is limited by the size of its context window. Methods such as text truncation, input segmentation, attention mechanism modifications, and memory summarization techniques have been proposed to address this limitation [46].

#### 2.2.1.2 Perception Module of LLM-based AI Agents

The perception module is the component of LLM-based AI agents that enables them to receive and process information from various sources and modalities [38]. This perceptual space allows the agents to understand their environment, and make informed decisions.

Textual input is the native modality that LLMs are highly proficient at processing. While LLMs inherently lack the ability to directly process visual information like images, integrating

visual perception significantly expands the agent's understanding of the world. One approach to enable LLMs to understand visual information is to combine an image encoder, such as those based on the vision transformer (ViT) [47] architecture, directly with the LLM to perceive visual content. The image is divided into patches, linearly projected, and treated as input tokens for the transformer. By calculating self-attention between tokens, information from the entire image can be integrated. While this end-to-end training of the visual encoder and LLM can achieve remarkable visual perception abilities, it comes at a substantial computational cost. Large language models that have been extended with these visual perception capabilities are referred to as vision large language models (VLLMs). A prominent example of a VLLM is GPT-4-Vision [48] from OpenAI, which incorporates an image encoder to gain the ability to natively process and understand visual information.

A more efficient paradigm is to use extensively pre-trained visual encoders and LLMs, freezing one or both of them during training to balance computational resources and model performance [49]. However, an extra learnable interface layer is typically required to align the visual encoder's output with the LLM's input embeddings. For example, the querying transformer (Q-Former) module has been used as an intermediate layer to extract language-informative visual representations.

Despite these advancements, current methods for visual-linguistic learning in LLMs tend to overfit to the pretraining dataset and struggle to generalize or few-shot learn in new domains [38]. This poses a particular challenge for niche applications like chart

understanding, where large-scale datasets for pretraining a visual encoder are not readily available [17]. The lack of a sufficiently diverse and comprehensive dataset of charts hinders the development of a robust visual perception module for this specific domain.

Given these limitations, an alternative approach to expanding the perception space of LLMs to the visual domain is by utilizing external tools with computer vision capabilities. This is made possible by the action module of the LLM-based agent, which enables the agent to interface with specialized visual processing tools to virtually extend its perception capabilities without the need for large-scale pretraining datasets or computationally intensive end-to-end training.

#### 2.2.1.3 Action Module of LLM-based AI Agents

While LLMs have demonstrated remarkable abilities in natural language understanding and generation, they often lack the domain expertise or computational capabilities to directly complete complex real-world tasks. One promising approach to address these limitations is enabling LLM-based agents to actively interact with their environment using external tools [38].

The ReAct framework, proposed by Yao et al. (2023) [50], is a methodology that enables LLMs to reason about and utilize specialized tools. It achieves this by combining reasoning and acting. ReAct augments the action space of an LLM agent to include both natural language generation and tool-specific actions. By interleaving reasoning traces and tool

interactions, the agent can dynamically compose information, track progress, and adapt its plans based on the evolving context. The thoughts guide the selection and application of tools, while the actions interface with the external tools to gather information or effect changes in the environment. Through this iterative process, the agent can break down complex tasks, reason about intermediate results, and construct robust and interpretable task-solving trajectories [51]. Due to its simplicity and impressive performance, the ReAct framework has been widely adopted in subsequent works [52, 46, 53].

Building upon the ReAct framework, Qin et al. (2023) [54] introduced ToolLLM, a framework for reinforcing open-source LLMs' capabilities in tool use through fine-tuning on a large-scale tool-augmented dataset. Central to the ToolLLM framework is the construction of ToolBench, an instruction-tuning dataset covering over 16,000 APIs. ToolBench characterizes each tool to the LLM by providing detailed documentation in a structured JSON format, including functionality descriptions, required parameters, and example responses. By pre-training on this rich tool dataset, LLMs can generalize their tool use capabilities to new, unseen tools through in-context learning, where the new tool is characterized in the same JSON format used during training. In addition to the fine-tuning framework, ToolLLM also proposes an alternative to the ReAct reasoning approach called depth-first search based decision tree (DFSDT) reasoning. In DFSDT, the agent explores multiple reasoning paths, retracts unproductive steps, and strategically navigates to the most promising path. While drastically more computationally intensive, this approach

enhances the planning and reasoning abilities of the LLM. While the ToolLLM framework has helped open-source models reduce the gap with closed-source models in terms of tool use capabilities, the latter remain significantly superior. In particular, GPT-4 currently achieves state-of-the-art results on the ToolBench benchmark, far outperforming open-source models like ToolLLaMA.

Recent work on Prismer [55] showcases how tool use can enable LLM-based agents to achieve state-of-the-art performance on vision-language reasoning tasks with orders of magnitude less data than prior approaches. By leveraging an ensemble of pre-trained vision, language, and multi-modal expert models as tools, Prismer efficiently pools their specialized knowledge to bootstrap strong multi-modal reasoning capabilities. The expert resampler and adaptor modules proposed in Prismer exemplify effective techniques for integrating multiple expert tools while preserving their pre-trained knowledge. Prismer's strong results demonstrate that tool use, when combined with judicious architecture design, can enable highly sample-efficient learning, robustness to noisy tools, and impressive zero-shot and few-shot generalization. These findings underscore the promise of tool use as a paradigm for modularizing multi-modal learning and facilitating knowledge transfer from pre-training to downstream tasks.

Importantly, the action module provides a means to extend an LLM-based agent's perception capabilities beyond its native language modality. By interacting with visual processing tools, the agent can indirectly perceive and reason about images without

requiring a specialized visual encoder to be integrated into its core architecture, offering a pragmatic alternative to the challenging and data-intensive process of directly training an LLM to encode images.

A prime example of such an LLM-based agent with extended visual perception through tool use is Visual ChatGPT [56]. This system integrates a diverse set of visual foundation models (VFMs) spanning vision and vision-language tasks with ChatGPT using a prompt manager module. The VFMs cover a wide range of visual processing capabilities, including visual question answering, image captioning, object detection, image generation, and image editing. The prompt manager, equivalent to the brain module in the framework from Xi et al. (2023) [38], serves as a bridge between the language and vision modalities, converting visual information into language instructions for selecting, executing, and chaining the VFMs based on the dialogue context. This enables Visual ChatGPT to perform tasks such as answering questions about image content, generating images from textual descriptions, and executing complex image editing operations.

Through multi-turn discussion experiments and case studies, Visual ChatGPT demonstrated impressive multimodal conversational and visual reasoning capabilities. It was able to handle multi-step visual editing instructions by sequencing relevant VFMs. For example, given a user request to "detect the depth map of the image first, show the depth image, then generate a new image based on this depth map, and finally describe the final image", Visual ChatGPT first used an depth VFM to compute the gradient-based depth of

the original image, an image-to-image VFM to generate from the depth image a new scene with the same character and objects, and finally a captioning VFM to describe the new scene of the image. In other experiments, the system also showed proficiency at answering questions about image content and style, such as identifying objects, describing their relationships, and making precise modifications to the image based on language input. The authors noted that the modular tool-based approach offers flexibility and extensibility advantages, allowing new VFMs to be easily incorporated to expand the agent's capabilities.

#### 2.2.2 Conclusion

In conclusion, LLM-based agents have emerged as a promising solution for autonomous tasks requiring cognitive capabilities, such as chart understanding and summarization. These agents operate through three specialized modules: a brain module for reasoning and decision-making, a perception module for processing various input modalities, and an action module for interfacing with external tools. Recent systems like Visual ChatGPT and Prismer have demonstrated the remarkable effectiveness of this approach.

# 2.3 Limits of Chart Summarization Evaluation Methods

# 2.3.1 Current Evaluation Methods in Chart Summarization and Their Limitations

Evaluation is central in validating research hypotheses and driving progress in any field. In chart summarization, evaluation aims to measure the quality of generated summaries, but faces significant challenges in deriving reliable, well-aligned metrics from this objective. This section examines current evaluation approaches in chart summarization, their limitations, and opportunities for improvement.

At its core, an evaluation procedure optimizes an objective - in this case, generating high-quality chart summaries. Abstractive text summarization aims to compress long textual documents into a short, human-readable form that contains the most important information from the source by detecting salient parts and paraphrasing them to form the final output, a concept that extends well to chart summarization [57]. However, measuring this objective poses difficulties. Effective metrics should be specific, aligning closely with the objective, and reliable, producing consistent results under identical conditions.

Evaluation methods for chart summarization fall into two main categories: quantitative heuristic scores and human evaluator pairwise comparisons. Heuristic metrics, such as

Bilingual Evaluation Understudy (BLEU), Recall-Oriented Understudy for Gisting Evaluation (ROUGE), and Consensus-based Image Description Evaluation (CIDEr), assess a candidate summary against a reference summary by comparing words or word groups [4]. While highly interpretable, these metrics heavily penalize paraphrasing and alternative word choices, limiting their flexibility. BERTScore [58] offers an alternative by comparing vector embeddings from a pre-trained BERT model, aiming to capture semantic similarity over exact wording. However, the opacity of neural networks renders this method less interpretable. All these metrics also require labeled reference summaries, which are scarce in chart summarization datasets. Consequently, heuristic scores often align poorly with the evaluation objective and show weak correlation with human judgments of summary quality [3, 59].

Human evaluator pairwise comparison, the second primary approach, is considered the gold standard in chart summarization evaluation [10, 4, 3]. Evaluators view a chart image and two candidate summaries, selecting their preferred summary. However, this method's high cost often constrains study size, failing to account for evaluator variability and chart diversity. Most studies, such as Chart-to-Text and UniChart, also do not specify the preference criteria given to evaluators, obscuring the basis for their judgments and compromising the method's specificity and objective alignment [4, 3].

The challenges in creating reliable and objective-aligned evaluation metrics for chart summarization are similar to those encountered in the more mature field of text

summarization. Text summarization has a longer history and greater resources, which have led to the development of strategies to address these common challenges. In the next section, we will explore how these strategies from text summarization evaluation could be adapted and applied to improve the evaluation of chart summaries.

# 2.3.2 Inspiration and Opportunities from Text Summarization Evaluation

Text summarization aims to compress long documents into shorter forms that capture the most important information. Two main approaches exist: extractive summarization, which selects and concatenates key sentences from the source text, and abstractive summarization, which generates new sentences that convey the essence of the original document [60]. As chart summarization more closely resembles abstractive summarization, this section focuses on evaluation methods in that domain.

Compared to chart summarization, text summarization benefits from a wealth of large-scale datasets. Notable examples include DUC 2004, containing 500 news articles paired with human-written summaries; CNN/DailyMail, featuring 300,000 articles with author-written highlight summaries; and XSum, a collection of 230,000 articles with single-sentence summaries. These large datasets have first pushed towards the development of automated evaluation metrics like ROUGE, BLEU, CIDEr, and BERTScore, which have subsequently been adopted in chart summarization. However, it was also noted in text summarization

literature that these metrics often poorly align with human judgments of summary quality [59].

To address these limitations, recent text summarization research has shifted towards more informative human evaluation protocols. While they used to only be pair-wise comparison, they now involve expert annotators grading specific quality criteria (QC) on a 5-point Likert scale, as proposed by, later reinforced in the SummEval meta-evaluation of metrics' quality [57, 59]. The four key criteria are:

- Coherence: Assessing the overall structure and logical flow of the summary.
- Fluency: Rating the grammatical correctness and readability of the summary.
- Consistency: Evaluating the factual alignment between the summary and source.
- Relevance: Determining if the summary captures the most important information from the source.

Introducing these evaluation criteria to chart summarization benchmarks has the potential to significantly enhance the informativeness and reliability of quality assessments. While this approach is considered the gold standard in text summarization, it is resource-intensive and currently difficult to automate, as highlighted by both Kryscinski et al. (2019) [57] and Fabbri et al. (2021) [59].

#### 2.3.3 Conclusion

In conclusion, current chart summary evaluation methods face significant limitations. Heuristic metrics like BLEU, ROUGE, and CIDEr often poorly align with human judgments, while human evaluator pairwise comparison lacks insight into specific quality attributes. Drawing from text summarization research, we propose addressing these challenges through quality criteria assessment by human annotators, enabling more reliable and informative quality measurements.

## Chapter 3

# Design of Large Language Model

# **Agent for Chart Summarization**

The field of chart summarization has seen significant advancements in recent years, primarily driven by end-to-end vision models such as UniChart. These models, while demonstrating promising results, rely heavily on extensive labeled datasets comprising chart-summary pairs. The creation and maintenance of such datasets present substantial challenges in terms of cost, time, and scalability, limiting the broader applicability of these approaches in real-world scenarios.

This chapter presents a paradigm for chart summarization that aims to address these limitations. Our proposed method leverages the expansive general knowledge, reasoning capabilities, and tool utilization proficiency of LLMs, combined with the precision and domain-specificity of chart visual question answering (CVQA) and chart derendering models. This approach forms an autonomous framework centered around an LLM "brain", capable of interpreting and summarizing charts without relying on expensive labeled datasets.

The primary objective of this chapter is to detail the design and architecture of a chart summarization system that circumvents the need for extensive labeled data. In doing so, we address the following research question: How can an LLM-based agent architecture effectively leverage domain-specific models as tools to enable chart summarization without labeled data? Our exploration focuses on a modular architecture comprising a brain module powered by an LLM and specialized tool modules, designed to enable flexible and efficient chart summarization.

The chapter is structured as follows: Section 3.1 provides an overview of the agent architecture and its key components. Section 3.2 presents the detailed implementation of each module, including the brain, action, and perception modules. Section 3.3 concludes with a discussion of the key design aspects, potential strengths and limitations of our approach, setting the stage for the empirical evaluation in the subsequent chapter.

### 3.1 Architecture Overview

The chart summarization system proposed in this thesis employs an autonomous agent type of architecture, inspired from works such as Prismer [55] and Visual ChatGPT [56], and

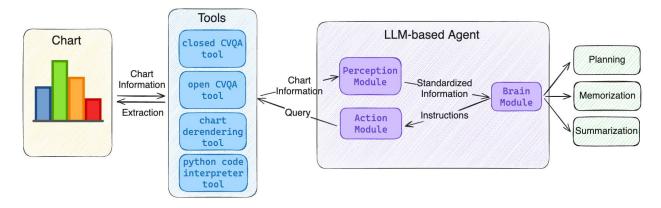


Figure 3.1: Architecture overview of the summarization system

following the guidelines established by Xi et al. (2023) [38], which combines the reasoning capabilities of LLMs with the precision of specialized visual processing tools. This approach aims to overcome the limitations of current end-to-end deep learning methods by leveraging a modular, tool-based design that enhances flexibility, improves sample efficiency, and increases interpretability. The architecture of the agent can be observed in Figure 3.1.

The system operates through an iterative process of information extraction, interpretation, and reasoning. The complete process follows this sequence: First, the brain module formulates plans and instructions based on its current understanding of the chart. Then, the action module executes these instructions by interacting with the appropriate tools. Next, the perception module receives and formats the extracted information, providing the brain module with formatted text for further analysis and decision-making. This cycle continues until the brain module determines that sufficient information has been gathered to generate a comprehensive chart summary, or until the predetermined iteration

limit is reached.

This modular design allows for the integration of various specialized tools, such as chart visual question answering (CVQA) models, chart derendering tools, and code interpreters.

These tools can be queried and utilized as needed by the action module under the guidance of the brain module.

The design choices for this architecture are motivated by several key rationales. Firstly, this modular tool-based approach offers significant theoretical advantages over traditional end-to-end deep learning methods. Firstly, it enhances flexibility and extensibility by allowing the system to adapt to new tools without the need for retraining. Tools are provided to the model in a few-shot learning style, accompanied by concrete explanations and examples. This adaptability enables the seamless integration of new machine learning models as they emerge, outperforming existing ones in specific domains of chart information retrieval such as label extraction, trend recognition, chart derendering, or CVQA. The section dedicated to the action module will present how and why we selected the specific tools used in our architecture.

Moreover, our approach has the potential to improve sample efficiency and accuracy by leveraging domain-specific deep learning models as tools, thereby reducing the demand for high-quality labeled data typically required to train end-to-end summarization deep learning models. Chart summarization is a complex task that demands strong capabilities in both visual information extraction and language understanding. Traditional approaches would

require an extensive dataset of chart-summary pairs, potentially in the order of millions, which is both costly and time-consuming to produce. In contrast, our agent framework could allow for the use of smaller, more specialized models trained on simpler tasks, such as CVQA and chart derendering, leaving the complex natural language summarization to the large language model. This division of labor potentially leads to improved sample efficiency, accuracy and reduced reliance on large-scale labeled datasets.

Finally, this approach can offer greater interpretability and controllability. Rather than relegating the reasoning and summary generation process to the abstract mechanisms of an end-to-end neural network, we leverage the agent's explicit reasoning capabilities. This allows for a more transparent identification of the information used to generate summaries and how it was interpreted to derive the high-level message conveyed by the chart. Such transparency facilitates easier error diagnosis, model improvement, and potentially, greater trust in the system's outputs.

The subsequent sections will dive deeper into each component of the architecture, delineating their functionalities, interactions, and the rationale behind their design choices.

## 3.2 Modules Implementation

#### 3.2.1 Perception Module

The perception module serves as an abstraction layer that translates diverse outputs from the action module's chart analysis tools into a format suitable for the brain module's processing. More concretely, in our implementation, the perception module is a set of functions which format tool outputs into a consistent string structure with predefined fields: a tool call ID, the tool's name, any arguments passed to the tool, and the tool's output which is prefixed with a standard "Observation:" marker. This structured information is then added to the LLM's prompt for the next iteration of brain reflection, ensuring consistent processing by the LLM while maintaining a traceable record of the agent's decision-making process through its logs

The capabilities of the perception module are inherently bounded by the underlying transformer architecture employed in our system. As most transformer-based LLMs during this research did not support direct image processing, our system processes chart images indirectly by converting visual information into text-based descriptions. Additionally, the transformer's context window limitation, typically ranging from 8,000 to 16,000 tokens [61], constrains the amount of information that can be processed simultaneously, leading us to implement a rolling summarization memory mechanism that will be detailed in the brain module section.

With this understanding of how our system processes and formats input information, we now turn to examine its core reasoning component: the brain module.

#### 3.2.2 Brain Module

The brain module uses GPT-4 (checkpoint 0613) as its core language model, selected for its superior performance in complex reasoning tasks and function calling capabilities [62].. This latter ability was essential for our tool-based chart analysis approach.

The second core component of the brain module is the system prompt that guides the model's reasoning process. This prompt instructs the LLM to analyze charts through alternating "Thought" and "Action" steps: reflection phases where it evaluates current information and identifies gaps, followed by action phases where it generates JSON commands for specific tools. Our Python backend parses and executes these commands and appends the results to the context for the next iteration. The effectiveness of this iterative process relies heavily on the prompt engineering techniques we examine next.

#### 3.2.2.1 Prompt Engineering Techniques

The brain module implements a combination of chain-of-thought (COT) [42], ReAct framework [50], and task decomposition as the main prompt engineering techniques. The system prompt structures the LLM's analysis process into explicit reasoning steps through task decomposition and iterative cycles of analysis.

When analyzing a chart, the brain module first breaks down the task of chart understanding into specific questions about visual elements, data points, and trends. Each sub-task then triggers a cycle of COT reasoning and tool interaction: the module reflects on what information is needed, selects appropriate tools to gather this information, and evaluates the results to inform its next steps. For instance, it might first reason about identifying the chart type and key elements, then systematically query tools about titles, labels, and specific data points, building a comprehensive understanding through each iteration."

The integration of these techniques in our prompt engineering serves multiple purposes. Task decomposition manages complexity by breaking down chart analysis into manageable components. COT reasoning ensures explicit and traceable decision-making in tool selection and information interpretation. The ReAct framework provides a structured way to interleave this reasoning with concrete actions, reducing potential hallucinations by grounding the LLM's analysis in tool-verified data. Central to this approach is the module's ability to effectively select and utilize appropriate tools throughout the analysis process, which we examine next.

#### 3.2.2.2 Tool Learning and Utilization

A major aspect of the brain module's functionality is its ability to effectively use a variety of tools for chart analysis. This capability is primarily achieved through ICL, a powerful feature of large language models that allows them to adapt to new tasks based on instructions and examples provided in their input prompt.

In our system, each tool is described to the brain module through a structured subprompt that includes a clear description of the tool's functionality, its input parameters, and the expected format of its output. We also provide several examples of the tool being used correctly in various scenarios, allowing the brain module to generalize from these examples to novel situations without requiring extensive training or fine-tuning. For instance, our "Python Code Interpreter" tool, which is useful for analyzing data from charts, might be described as follows:

"The 'PythonCodeInterpreter' tool executes Python code on a pandas dataframe representation of the chart data. It has access to numpy, pandas, and scipy libraries, as well as a pre-defined 'df' dataframe containing the chart data, obtained from the 'ChartDerendering' tool. A custom 'get\_slope' function is available for trend analysis. The tool takes Python code as input and returns the printed output. Example usage: Input: print(df['Value'].mean(), get\_slope(df, 'Value', window\_size=2)), Output: 42.5, 0.75."

This approach is particularly valuable for tools like the Python Code Interpreter, which need to perform a wide range of data analysis tasks, from simple statistical calculations to complex trend analyses, adapting their approach based on the specific characteristics of each chart.

#### 3.2.2.3 Memory Mechanism

As mentioned in the section about the perception module, the brain module of our LLM-based agent faces a significant challenge in the form of context length limitations. The GPT-4 checkpoint (0613) employed in our system has a maximum context length of approximately 8,000 tokens (about 6,000 words). This constraint poses a potential issue for our chart summarization task, which often requires multiple rounds of reasoning and tool interactions. As the analysis progresses, there's a risk of "forgetting" valuable information from earlier rounds, potentially compromising the quality and coherence of the final summary.

To address this challenge, we explored several memory management strategies. Initially, we considered selective information retention, which involved filtering the types of information kept in memory. However, our experiments revealed that preserving both reasoning traces and tool calls was central for maintaining coherence between reasoning rounds and enhancing the agent's overall performance. We also explored information compression by using an LLM to summarize reasoning traces, but this approach led to a slight decline in inter-round coherence and significantly increased processing time due to additional LLM calls.

Based on our experimental results, we adopted a rolling window with summarization strategy. This approach offers a balance between memory conservation and information retention. As the agent progresses through reasoning cycles, it accumulates information in its working memory. When the accumulated information approaches the context length

limit (approximately 7,000 tokens in our implementation), the oldest information is summarized into a concise yet insightful paragraph. This summary replaces the original detailed information, freeing up memory space while preserving key insights. The process repeats as necessary throughout the chart analysis task.

In practice, our agent rarely encounters situations where it exceeds the context window, as most chart analyses can be completed within the available token limit. Still, for particularly complex or data-rich charts that do approach this limit, this approach proves effective. While this method works well in maintaining coherent reasoning, it inevitably results in some loss of detail through the summarization process.

It is worth noting, however, that newer language models have been developed with significantly longer context windows [61], largely mitigating the need for such complex memorization strategies.

#### 3.2.2.4 Summary Generation

The summary generation phase is implemented through a prompt distinct from the information extraction phase. This prompt receives all previously collected information; including tool calls, their inputs and outputs, and reasoning traces; and instructs the LLM to generate a coherent summary. The effectiveness of this phase heavily relies on careful prompt engineering, which we detail in this section.

Using COT reasoning, the LLM analyzes the compiled information to plan and structure

the summary, ensuring systematic integration of all relevant details. Our prompting strategy builds on Kintsch and van Dijk's model of text comprehension [63], which emphasizes the importance of information organization for effective understanding. The prompt provides specific instructions for the summary structure: an overview of the chart's main topic, key trends or data points, and notable implications. This structured approach ensures both consistency across chart types and deep understanding of the chart's message.

The prompt also includes explicit quality criteria to guide the generation of the LLM. These criteria encompass coherence, consistency, fluency, and relevance. By explicitly stating these criteria, we guide the LLM towards producing summaries that not only convey factual information but also present it in a clear, well-structured, and meaningful manner.

To further enhance quality and consistency across various chart types, we leverage fewshot learning [41] by incorporating three manually crafted and carefully reviewed examples of high-quality summaries into our prompt. These examples serve as implicit templates, demonstrating the practical application of our structural approach and quality criteria for different chart types and complexities.

Drawing on Miller's seminal work on cognitive limitations [64], we implement a target word count of 50-70 words per summary in our generation process. This constraint serves multiple purposes: it ensures conciseness while aligning with human cognitive processing capacity; it compels the system to prioritize the most salient information; and it maintains consistency across diverse chart types, facilitating user familiarity and expectations.

In situations where the brain module cannot discern a clear overarching message from the extracted information, such as when the data points are highly uncorrelated or contradictory, we employ a specialized prompting strategy. The LLM is instructed to analyze the available data and generate a summary that explicitly acknowledges the uncertainty of its conclusions. For instance, the summary might begin with phrases like "The chart presents varied data points without a clear trend..." or "While the information is inconclusive, possible interpretations include...". This approach mitigates the risk of hallucination by encouraging the model to express uncertainty rather than fabricate connections. Moreover, by explicitly stating uncertainty in certain situations, we implement a form of negative framing that, according to MacLeod et al. [65], encourages appropriate skepticism in BVIP towards AI-generated content.

After the initial summary generation, we implement a self-reflection mechanism inspired by recent research on AI self-improvement [46]. This approach relies on the LLM's capability to critically evaluate its work against the established quality criteria. Based on this self-assessment, the LLM revises and regenerates the summary. Our experimental tests have shown improvement after the first iteration of this process, but diminishing returns thereafter, leading us to limit the process to a single iteration that balances quality with computational efficiency.

While this approach aims to produce informative and accessible summaries, the quality of the initially extracted information significantly impacts the model's ability to generate insightful content. The action module, which executes the brain's instructions and interfaces with external tools, is therefore crucial to the overall summarization process.

#### 3.2.3 Action Module

The action module serves as the third core component of our LLM-based agent architecture, acting as the interface between the brain module and external tools for information extraction. It extends the agent's perception capabilities by managing tool interactions: handling API requests, error management, load balancing, and parallel execution of tool operations.

During system initialization, the action module provides the brain module with information about available tools, including their parameters and usage instructions. This modular design allows for easy integration of new tools. This section details the implementation of the action module and provides a comprehensive overview of the integrated tools.

#### 3.2.3.1 Tool Interaction Protocol

The action module implements a standardized JSON-based protocol for communication between the brain and action modules. The protocol defines a simple structure: a tool identifier and a parameters dictionary containing the necessary information for the tool's function. For instance, the brain module might provide the following execution instruction

for the UniChart ClosedCQA tool:

```
{
   "tool_id": "unichart_closedcqa",
   "parameters": {
        "question": "What is the highest value in the chart?",
        "chart_image": [image_data]
   }
}
```

The action module would then process this instruction, make the appropriate call to the UniChart ClosedCQA tool, and return the result to the brain module. The output might look like this:

```
{
    "status": true,
    "answer": "125 million sales",
    "error_message": null
}
```

For reliable operation, we implemented error handling in the action module. It first validates all inputs before execution, ensuring required parameters are present and correctly formatted. For potential failures, the system implements a retry mechanism with increasing delays between attempts, allowing up to three retries within a 30-second window to prevent tool overload. The module maintains detailed logs of all interactions, facilitating system improvement and troubleshooting.

The action module implements several optimization techniques to enhance performance. When the brain module requests multiple independent tool executions, these are processed in parallel to reduce overall execution time. The system also implements a caching mechanism that stores summaries alongside chart hash signatures, enabling quick retrieval of previously analyzed charts. Additionally, multiple independent calls to the same tool are batched and executed together for improved efficiency.

#### 3.2.3.2 Tools List

The action module integrates several domain-specific tools for chart analysis. We selected the tools based on three factors: quality of extracted information, coverage of different information types, and ability to cross-verify information through complementary capabilities.

The UniChart model forms our toolset's backbone, providing multiple modes: CVQA, chart derendering, and chart summarization. These modes share the same self-supervised pretraining checkpoint, obtained through standard encoder-decoder pretraining [47], but are fine-tuned for specific tasks. We utilize the CVQA and chart derendering modes to decompose complex chart summarization into manageable sub-tasks. This approach achieves sample efficiency in two ways: the pretraining leverages readily available unlabeled chart images from the internet, while the fine-tuning only utilizes established CVQA and chart derendering labeled datasets, thus avoiding the need for expensive chart summary

annotations.

The UniChart ClosedCVQA mode complements the OpenCVQA mode by answering close-ended questions about specific chart values and attributes. It's capable of extracting definite facts and provides high-confidence responses for quantitative information. One of its limitations is that, when identifying the highest value, it might return only one instance even if multiple exist. It may also struggle with questions requiring contextual understanding beyond explicit data.

The UniChart Chart Derendering mode converts visual chart data into structured, tabular format. To enhance its capabilities, we implemented a post-processing step using an LLM that cleans the table and creates consistent formatting across chart types, addressing output variability in the original model. Empty values are cleaned and appropriate value conversions are performed for compatibility with Python and Pandas libraries. Since data table extraction is typically required for each chart, we precompute it during system initialization, allowing immediate access through the Python Pandas Interpreter tool. While this mode may struggle with unusual chart designs and dense data points, it provides essential raw data accessibility.

The Python Pandas Interpreter tool executes custom Python code on the dataframe generated by chart derendering. It utilizes Pandas and NumPy libraries for advanced statistics and derived insights. We implemented a custom function that fits a linear model to two series and returns the fitting score and proportionality coefficient, addressing a common need in chart interpretation. Safe code execution is ensured through a combination of prompt engineering and a restricted Docker environment that runs without privileged access, filesystem writes, or network connectivity.

All tools are implemented as Python applications with HTTP endpoints, containerized using Docker for standardized interaction and reproducibility. The models are served on a single GPU with 8GB of VRAM.

#### 3.3 Discussion

#### 3.3.1 Design Choices and Trade-offs

Our modular LLM-based agent architecture addresses a fundamental challenge in chart summarization: the scarcity of high-quality chart-summary pairs. Rather than relying on extensive labeled datasets, we leverage the general knowledge and reasoning capabilities of LLMs combined with specialized chart analysis tools. This design choice achieves sample efficiency through strategic decomposition of the summarization task: utilizing pre-trained LLMs for reasoning and specialized tools for specific chart analysis tasks avoids the need for end-to-end training on chart-summary pairs. The efficiency extends to practical deployment, where new chart types or domains can be handled by updating individual tools or prompts rather than retraining an entire system.

When compared to end-to-end approaches like UniChart or VLLMs like GPT-4V, our

system offers distinct trade-offs. While potentially more complex in terms of system integration, this complexity is balanced by greater flexibility and interpretability. Each component can be independently improved or replaced, allowing the system to evolve with technological advances. End-to-end models, though providing seamless integration between visual understanding and language generation, require large amounts of task-specific training data and offer limited flexibility for improvement. VLLMs, while powerful, demand even more extensive datasets and typically remain inaccessible for open-source development due to their size and computational requirements.

The primary technical challenge lies in coordinating multiple specialized tools while ensuring consistent information flow. Each tool introduces its own limitations and potential biases, which must be carefully managed through robust error handling and cross-verification mechanisms. The effectiveness of the system heavily relies on the brain module's ability to select and utilize appropriate tools for each chart type. However, LLMs have demonstrated strong capability in reconciling potentially inconsistent information from various sources, making them well-suited for this coordination role.

The use of AI for chart summarization, particularly for BVIP users, raises important ethical considerations that influenced our design choices. Potential biases in LLMs or domain-specific tools could disproportionately affect BVIP users who rely entirely on the generated summaries. We address this through specific mechanisms: explicit uncertainty communication in unclear cases, verification of extracted information through multiple

tools, and careful prompt engineering to encourage balanced and accurate reporting. This transparency is particularly valuable in applications for BVIP users, where understanding the system's limitations and behavior is crucial for building trust and ensuring reliable usage.

#### 3.3.2 Future Directions and Implications

Our modular approach demonstrates how complex AI systems can be designed without relying on extensive task-specific datasets, with broader implications for developing AI accessibility tools where labeled data is often scarce but the need for reliable and transparent solutions is crucial.

Several promising research directions emerge from this work. The integration of emerging multimodal LLMs could enhance direct visual understanding while maintaining our system's flexibility. Specialized tools could be developed for specific chart types or domains, such as scientific publications or financial reports. The prompt engineering techniques we developed could be formalized into reusable patterns for similar modular AI systems. Additionally, investigating methods for automated tool selection and execution optimization could further improve the system's efficiency.

While our architectural analysis suggests potential advantages, empirical validation is essential. The next chapter evaluates our system's performance, focusing particularly on summary quality and accessibility for BVIP users.

## Chapter 4

# System Evaluation

Building upon the foundation laid in Chapter 3, where we introduced a novel LLM-based agent architecture for chart summarization, this chapter aims to provide a thorough evaluation of our proposed system against relevant baselines, namely UniChart and GPT-4-Vision. The evaluation process is designed not only to gauge the performance of our system but also to contribute to the broader discourse on effective assessment methodologies in the field of chart summarization.

Traditional evaluation approaches in chart summarization face significant limitations. Automated metrics like BLEU and ROUGE, while widely used, rely on expensive labeled datasets and often fail to capture the subtle aspects of summary quality, particularly in visual-textual contexts. Human evaluator rankings, though valuable, are resource-intensive, limited in scale, and provide little insight into the reasons behind preferences. To address

these challenges, we adapt quality criteria from text summarization to provide a more informative and multifaceted assessment of chart summaries. This adaptation allows us to precisely characterize chart summaries and investigate the key criteria influencing user satisfaction across different user groups.

In this chapter, we aim to address two research questions that are central to our evaluation of chart summarization systems and methodologies. First, we seek to determine how our modular LLM-based agent approach compares to end-to-end visual-language training and VLLMs in terms of summarization quality. Second, we seek to explore the key factors that influence user satisfaction across different groups, as assessed by our adapted quality criteria.

The remainder of this chapter provides an evaluation of our chart summarization system, addressing the aforementioned research questions through a structured approach. We begin by examining the dataset used for evaluation in Section 4.1, focusing on its characteristics and relevance to real-world scenarios. Section 4.2 then outlines the baseline methods selected for comparison, justifying their inclusion based on performance and relevance. Our evaluation methodology is detailed in Sections 4.3 and 4.4. The former introduces the quality criteria framework adapted for this study, explaining each criterion and its application. The latter describes the design and implementation of our user study, encompassing participant selection, materials preparation, and data collection procedures. Section 4.5 presents an analysis of the evaluation results, examining both system performance and user preferences. The chapter concludes with Section 4.6, which discusses

the implications of our findings, acknowledges limitations, and proposes directions for future research in chart summarization evaluation.

### 4.1 Dataset

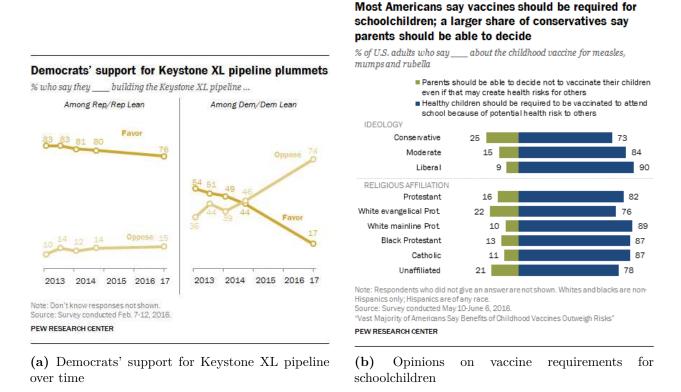


Figure 4.1: Examples of Charts from Pew Chart2Text Subset

For the evaluation of our chart summarization system, we utilized the Pew subset of the Chart-to-Text dataset [4]. This choice was motivated by several factors that align with our research objectives and the real-world applicability of our proposed LLM-based AI agent approach.

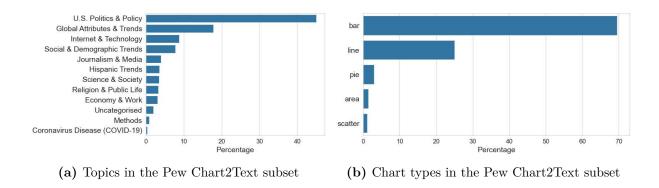


Figure 4.2: Distributions of topics and chart types in the Pew subset of Chart2Text dataset

The Pew subset comprises 10,000 real-world charts extracted from articles published by the Pew Research Center. For our study, we focused specifically on the 1,200 charts from the test split, ensuring a robust evaluation of our system on unseen data. This dataset offers a diverse and representative sample of charts encountered in actual research and journalism contexts, making it particularly suitable for assessing the performance of our chart summarization system in practical scenarios.

One of the key strengths of this dataset is its topical diversity. The charts span 12 distinct categories, covering a wide range of subjects typically addressed in social and economic research. This variety is crucial for evaluating our system's ability to generate accurate and contextually appropriate summaries across different domains, a capability that is essential for a versatile chart summarization tool. In terms of chart types, the dataset primarily consists of bar and line plots, which together account for over 90% of the total number of charts. Additionally, it includes pie charts, area plots, and scatter plots, albeit in smaller proportions. This distribution reflects the prevalence of different chart types in real-world publications

and allows us to assess our system's performance across various visual representations of data.

Each chart in the dataset is accompanied by rich metadata, including the chart title, chart type classification, bounding boxes around main visual elements, and the underlying data table of values represented in the chart. While the dataset also includes captions for each chart, previous research has identified significant limitations in their quality. These captions, generated through automated heuristics that selected relevant surrounding paragraphs, often miss important information or include details not present in the chart itself, potentially leading to serious hallucination problems for models trained on these captions [4]. Given these limitations, we have made a deliberate decision to exclude these captions from our evaluation process. Instead, our research will focus exclusively on the chart images and the aforementioned metadata annotations. This approach allows us to leverage the most reliable components of the dataset, ensuring that our evaluation is based on accurate information.

The selection of this dataset and its components allows us to assess our chart summarization system's performance in a context that closely mirrors real-world scenarios. With this robust dataset in place, we can now turn our attention to establishing appropriate evaluation baselines, which will serve as crucial reference points for measuring the effectiveness of our proposed approach.

### 4.2 Baselines

In order to evaluate the effectiveness of our proposed chart summarization system, we have selected two state-of-the-art models as baselines for comparison: UniChart and GPT-4V.

UniChart [3] is an end-to-end visual-language model specifically designed for chart understanding and summarization. It was trained on a large, diverse dataset of over 627,000 charts using a multi-stage learning approach. The model's architecture consists of a visual encoder to process chart images and a text decoder to generate outputs. UniChart's training procedure involved an initial pretraining phase on multiple chart-related tasks, including data table generation, numerical and visual reasoning, open-ended question answering, and chart summarization. This comprehensive pretraining strategy aimed to imbue the model with a broad understanding of chart elements and their relationships. Following pretraining, UniChart was fine-tuned to create three specialized models, each focusing on a different aspect of chart understanding: chart visual question answering, chart derendering, and chart summarization. These models were fine-tuned on a public benchmark dataset, specifically the Chart-to-Text corpus, which includes charts from sources like Statista and Pew Research Center. This fine-tuning phase allowed each model to specialize in its respective task while building upon the foundation established during pretraining. This model held the state-of-the-art position in chart summarization until late 2023.

We chose UniChart as our first baseline due to its direct relevance to the chart

summarization task and its impressive performance. However, it is important to note that UniChart's reliance on a large, labeled dataset of chart-summary pairs highlights one of the key motivations for our research: to develop a system that can perform well without the need for extensive, task-specific labeled data.

The second baseline we selected is GPT-4V [48], a large-scale pretrained VLLM by OpenAI. Released in late 2023, GPT-4V extends the capabilities of the GPT-4 language model to include visual understanding. While the exact details of its training procedure are not publicly available, it is reasonable to assume that GPT-4V was trained on a vast corpus of text and images, likely largely surpassing the scale of UniChart's training data.

GPT-4V was selected as a baseline due to its state-of-the-art performance in visual-language tasks and its potential suitability for chart summarization. Given the scale of its pretraining, it likely encountered numerous chart images, and considering the base GPT-4 model's proficiency in text summarization, we can reasonably infer GPT-4V's capability in chart summarization. This combination of visual understanding and presumed summarization skills makes GPT-4V a compelling baseline, despite the lack of transparency regarding its specific training on chart-related tasks.

By comparing our proposed system against both UniChart and GPT-4V, we aim to evaluate its performance from different perspectives. UniChart provides a comparison against a specialized, end-to-end supervised model, while GPT-4V offers a benchmark for performance against a general-purpose, large-scale VLLM. This dual comparison will help

us assess the effectiveness of our approach in leveraging AI agents with domain-specific tools for chart summarization.

While these baselines offer strong points of comparison, effectively evaluating chart summarization models requires a tailored approach. Traditional text summarization metrics may not fully capture the nuances of chart summarization, where visual elements play a crucial role. To address this, we adapt an existing quality criteria evaluation framework from the field of text summarization to the context of charts. In the following section, we present this adapted framework.

## 4.3 Quality Criteria Evaluation Framework

The evaluation of chart summarization presents unique challenges that require a specialized framework. While text summarization has established evaluation protocols, chart summarization methods are comparatively underdeveloped. This section presents a quality criteria evaluation framework for assessing chart summaries.

The motivation for adapting text summarization criteria to chart summarization stems from several factors. Traditional metrics such as BLEU and ROUGE have shown poor correlation with human judgment of summary quality in the context of charts [59]. These metrics often fail to capture details that significantly impact summary quality [66]. Moreover, their reliance on labeled data, which is scarce and expensive to obtain for chart summarization, limits their scalability and applicability in this domain.

While pairwise comparison of summaries by human evaluators is considered the gold standard in chart summarization evaluation due to its reliability, it lacks granularity in explaining preferences and characterizing summary styles. To address these limitations, we propose adapting the quality criteria framework from text summarization, as defined by Kryscinski et al. (2019) [57], to the context of chart summarization.

The framework comprises four key criteria: coherence, fluency, consistency, and relevance.

In the context of chart summarization, these criteria are defined as follows:

- Coherence: Collective quality of the sentences made the overall structure and logical flow of the summary.
- Fluency: Quality of individual sentences as in the grammatical correctness and readability of the summary.
- Consistency: Factual alignment between the summary and the source.
- Relevance: Selection of the most important information from the source.

To adapt this framework for chart summarization, we simply redefine the "source" as the chart image, rather than the text documents typically used in traditional summarization tasks.

Fabbri et al. (2021) [59] established a protocol for applying the quality criteria framework in text summarization evaluation. Their process began with the random selection of 100 articles from the CNN/DailyMail test set, a widely used benchmark dataset for text summarization consisting of news articles paired with human-written

summaries, ensuring a diverse sample. For evaluator recruitment, they adopted a dual approach: engaging crowd-sourced annotators with professional English proficiency and a proven annotation track record, while also involving three expert annotators with academic This strategy balanced scalability with expertise. experience in summarization. The evaluation conducted online platform where was on an annotators rated computer-generated summaries on a 5-point Likert scale for each quality criterion, with each summary assessed by eight different evaluators to ensure reliability. Their analysis utilized Krippendorff's alpha coefficient to measure inter-annotator agreement and computed correlation scores between human ratings and automated metrics. This methodology revealed that while most automated metrics correlated poorly with human judgment, the proposed criteria showed high correlation with human preference, were interpretable, and measured distinct aspects of summary quality. By averaging scores for each criterion, Fabbri et al. compared model performance across these dimensions, demonstrating the framework's utility for evaluation purposes.

The quality criteria framework, adapted from text summarization, provides a structured approach to evaluating chart summaries. The following section describes the design of our user study, which employs this framework to compare the performance of chart summarization models.

## 4.4 User Study

To evaluate the effectiveness of our proposed LLM-based AI agent approach for chart summarization, we conducted a user study with three primary objectives. First, we aimed to collect quantitative data on the quality of model outputs using the adapted quality criteria framework described in Section 4.3. This assessment would provide insights into how well each model performs across the dimensions of coherence, fluency, consistency, and relevance. Second, we sought to gather data on user preferences by having participants rank summaries from different models. By comparing these rankings with the quality criteria ratings, we could identify which factors most strongly influence user satisfaction. Finally, we collected demographic and expertise data to investigate whether domain knowledge impacts user preferences and quality assessments of chart summaries.

Ethical considerations were important in the design and execution of our study. We obtained approval from the university's research ethics board (REB) prior to commencing the study. Participants were provided with clear information about the study's purpose, duration, and the nature of their involvement, and were offered a compensation of \$15 for their participation. Informed consent was obtained from all participants, and they were assured of their right to withdraw at any time without penalty. To protect participant privacy, all data collected was anonymized and stored securely. Participants were informed that the study carried risks similar to those involved in participating in any video conference or regular computer use, including potential discomfort when discussing feelings related to

technology use.

### 4.4.1 Participant Selection and Categorization

An important part of our study design was the comparison of responses between expert and novice users in the domain of the charts being summarized. This focus was motivated by previous research in text summarization, which suggests that experts and novices may have differing preferences and needs when it comes to summary content and style [67]. We hypothesized that similar differences might exist in the context of chart summarization. Understanding these potential differences is valuable for the development of summarization systems that can effectively serve diverse user groups.

Given the range of topics covered in the Pew Research Center Dataset, we selected "Economy & Work" as our primary domain of study. This choice was driven by several factors. First, the field of economics often requires specialized terminology, understanding of complex relationships between multiple variables, and substantial background knowledge, making it an ideal domain to distinguish between expert and novice users. Second, this topic offered sufficient diversity in chart types, allowing us to evaluate our system's performance across various visual representations of data. Lastly, economic data is often complex and intricate, presenting a challenging test case for our summarization system.

To recruit participants with varying levels of expertise, we targeted different university departments. For potential experts, we reached out to the departments of Economics, Management in Data Analytics, Finance, and Politics. To recruit novices, we approached faculties less likely to have extensive economics knowledge, such as the Faculty of Arts in which subjects of literature, art, and philosophy are taught and studied. This strategy aimed to ensure a diverse pool of participants with a range of domain knowledge.

To objectively classify participants as experts or novices, we developed a knowledge questionnaire comprising 15 multiple-choice questions of varying difficulty levels: four easy, seven medium, and four hard. For each question, participants had to choose one answer among four possibilities. The questionnaire was designed to ensure its validity in assessing economic knowledge. To discourage random guessing on harder questions, which reward more points, we included an "I don't know" option for the hard questions only. We did not offer this option for medium and easy questions to encourage participants to think through these questions rather than immediately selecting "I don't know" when unsure. Correct answers were awarded points based on their difficulty: one point for easy questions, two points for medium questions, and three points for hard questions. The questionnaire's difficulty levels were verified and calibrated with input from a doctoral student in the Economics department. The expertise score was calculated using the following formula:

Expertise Score = 
$$\frac{e \times 1 + m \times 2 + h \times 3}{4 + 7 \times 2 + 4 \times 3} \times 10$$

Where e, m, and h represent the number of correct answers in the easy, medium, and hard categories, respectively. This scoring system yields a range from zero (minimal expertise) to

ten (maximum expertise). We classified participants as novices if they scored four or below, experts if they scored seven or higher, and intermediate for scores between four and seven.

### 4.4.2 Materials Preparation

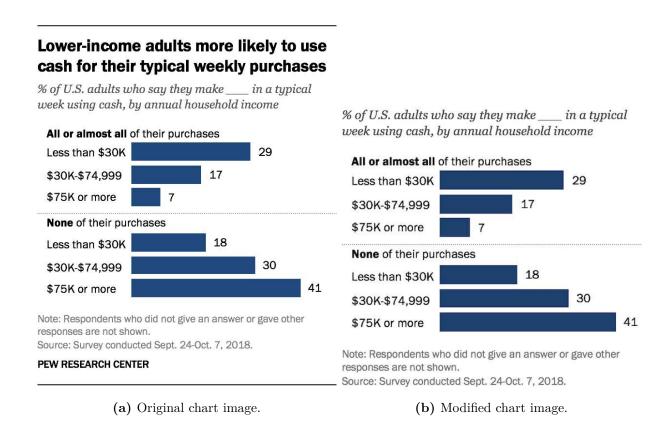


Figure 4.3: Comparison of original and modified charts.

The preparation of materials for this study involved three primary components: the selection and modification of charts, the generation of chart summaries using different systems, and the compilation of these elements into the evaluation form.

For this study, we selected ten charts from the Pew Research Center test set, specifically

from the "Economy & Work" category. The selection comprised four bar charts and six line charts, presenting a variety of styles and content. This mix was intended to provide a comprehensive assessment of the summarization models' capabilities across different chart types commonly used in economic reporting.

The decision to limit the study to ten charts was based on considerations of participant fatigue and time constraints. Our preliminary assessments indicated that evaluating and ranking three summaries for each chart would be a cognitively demanding task. We determined that ten charts would strike an optimal balance between gathering sufficient data and maintaining the quality of participant responses, while keeping the total study completion time to approximately one hour per participant.

Each selected chart underwent minimal modifications. We only cropped the original titles from the charts, as these often explicitly stated the main message of the visualization. This decision was motivated by two factors: firstly, it more closely simulated real-world scenarios where charts may not always have clear, descriptive titles; secondly, it presented a more challenging task for the summarization systems, requiring them to derive the main message solely from the visual elements and data presented in the chart. Figure 4.3 illustrates this preprocessing step, showing an example of an original chart and its cropped version.

Following the chart preparation, we generated summaries using the three systems described in Section 4.2: UniChart, GPT-4-Vision, and our proposed LLM-based summarization agent. The process for each model was as follows:

- UniChart: This neural network model received the cropped charts as direct input and generated summaries based on its training.
- GPT-4-Vision: We provided this model with the same prompt designed for the summary generation phase of our agent (as detailed in Section 3.2.2.4).
- LLM-based Summarization Agent: Our proposed model processed the cropped charts following the multi-step approach described in Chapter 3, which involves information extraction using domain-specific tools followed by summary generation using a LLM.

The prepared charts and their corresponding summaries were then assembled into an online Microsoft Form. This form was structured in two main parts to facilitate a comprehensive evaluation process. The first part began with a brief demographic questionnaire, designed to assess participants' familiarity with data visualizations and their primary sources of knowledge about economic topics. The expertise assessment quiz, as described in Section 4.4.1, was also integrated into this introductory section of the form.

The second part of the form focused on the evaluation of chart summaries. This section presented the chart images alongside their corresponding summaries for preference ranking and quality criteria rating. The specific layout and methodology for this evaluation process will be elaborated in the following section on data collection procedures.

#### 4.4.3 Data Collection Procedures

Prior to the evaluation, participants received a detailed email outlining the study procedures. This communication served to prepare them for the task and ensure consistent understanding across all evaluators. The email informed participants about the study's purpose of improving access to visual data for visually impaired individuals through automated chart summarization. It also set clear time expectations, advising that the evaluation would take approximately one hour, with 10-15 minutes allocated for the initial questionnaire and 45-50 minutes for the main evaluation of chart summaries.

An important component of the preparation was the introduction of the four quality criteria: coherence, fluency, consistency, and relevance. Each criterion was presented with its definition and accompanying guiding questions for participants to consider during their evaluations. To further enhance understanding of the grading process, participants were provided with a sample chart image and examples of poor summaries for each quality criterion. This approach was chosen to calibrate participants' expectations and ensure more consistent evaluations across the study. The email also provided instructions for participants to rank the summaries based on their personal preference.

To minimize external variables and ensure focus, participants were asked to complete the evaluation in a quiet setting and in one sitting. Contact details for technical or procedural support were provided, along with information about the \$15 compensation upon completion, acknowledging participants' time and effort.

The evaluation process was structured to collect both quantitative ratings and qualitative feedback for each chart summary, while minimizing potential biases. To avoid order effects, all charts were presented to participants in a randomized sequence. For each chart, participants were first shown the image and instructed to thoroughly understand its main message before proceeding. This step ensured that evaluators had a solid grasp of the chart's content, enabling more accurate assessments of the summaries.

Following this, all summaries for a given chart were presented simultaneously in a randomized order, without identifying which model generated each summary. This blind presentation method was chosen to minimize bias and encourage objective comparisons between summaries. To ensure consistent evaluation across all charts and summaries, participants were reminded of the quality criteria definitions before each assessment.

The quantitative assessment consisted of two components. First, participants rated each summary on a 5-point Likert scale (1 = poor, 5 = excellent) for each of the four quality criteria. In addition, after completing the individual ratings, participants ranked the summaries in order of preference. To capture the reasoning behind these rankings, participants were then invited to optionally leave a comment explaining their thought process or sharing any observations about the summaries themselves.

This process was repeated for all charts in the study. Upon completion of the entire evaluation, participants were thanked for their time and valuable input. To facilitate a more in-depth exploration of participants' experiences, they were asked if they would be willing

to participate in a follow-up interview. However, none of the target participants, expressed willingness to engage in these interviews, so they did not take place. Finally, participants were given the opportunity to leave a comment about their overall experience with the study or share any additional thoughts they had.

Section 4.5 presents an analysis of the collected data. This analysis encompasses the performance evaluation of different models across the established quality criteria and an investigation of user preferences.

### 4.5 Results

This section presents the outcomes of our user study on chart summarization, focusing on the performance of our proposed LLM Agent compared to two baseline models: GPT-4-Vision and UniChart.

Our study involved twenty participants: two novices, twelve intermediates, and six experts in fields related to economics and data analysis, as determined by the expertise questionnaire. The slight skew towards higher expertise levels can be attributed to the fact that many respondents were graduates from fields targeted for expert recruitment. The gender distribution was relatively balanced, with eleven men and nine women. Each participant evaluated summaries for ten different charts, each generated by three distinct models. This process yielded a total of 600 individual summary evaluations, providing a comprehensive dataset for analysis.

In the following subsections, we examine the user preferences (Section 4.5.1) and quality criteria (Section 4.5.2) in detail. These analyses offer insights into model performance across various dimensions of summary quality and user satisfaction, considering the impact of participant expertise levels on these assessments.

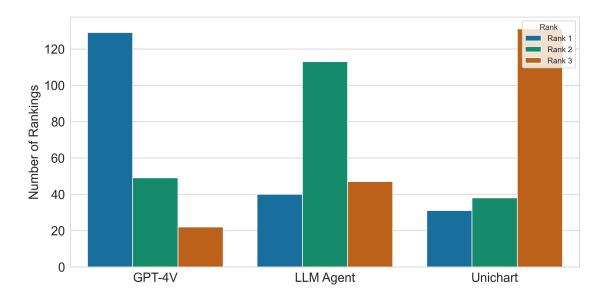
### 4.5.1 User Preferences

As part of our comprehensive evaluation, we analyzed user preferences to gain insights into how participants ranked the different chart summarization systems.

Figure 4.4 presents the distribution of rankings for each model across all charts and participants. The bar chart illustrates the number of times each model (GPT-4V, LLM Agent, and UniChart) was ranked first, second, or third in the evaluation process.

As observed on the figure, GPT-4V consistently received the highest number of first-place rankings, indicating a strong overall preference among participants. LLM Agent predominantly received second-place rankings, while UniChart was most frequently ranked third.

To quantify the level of agreement among raters and assess the statistical significance of these preferences, we employed several statistical tests. First, we calculated Kendall's W coefficient to measure the degree of agreement among raters across all scenarios. The average Kendall's W value was 0.3317, indicating a moderate level of agreement among participants. This suggests that while there is some consensus in the rankings, there is also a degree of



**Figure 4.4:** Ranking Distribution of Chart Summarization Models Across Evaluation Criteria. This figure illustrates how often GPT-4V, LLM Agent, and UniChart achieved each rank (1st, 2nd, or 3rd) in the evaluation. GPT-4V most frequently achieves the first rank, the LLM Agent predominantly occupies the second rank, and UniChart is most often placed third.

variability in individual preferences.

To determine whether the observed differences in rankings were statistically significant, we conducted a Friedman test. The test yielded a statistic of 107.49 with a p-value of  $10^{-24}$ , which is well below the conventional significance level of 0.05. This result provides strong evidence of statistically significant differences in the rankings of the three models.

Given the significant result of the Friedman test, we proceeded with a post-hoc analysis to identify specific pairwise differences between models. We employed the Nemenyi post-hoc test, which is specifically designed for use after a Friedman test and provides a conservative

Comparison	p-value	Significant
GPT-4V vs. LLM Agent	$10^{-4}$	Yes
GPT-4V vs. UniChart	$10^{-4}$	Yes
LLM Agent vs. UniChart	$10^{-4}$	Yes

Table 4.1: Nemenyi Post-hoc Test Results

approach to control for family-wise error rate in multiple comparisons. The results of the Nemenyi test are presented in Table 4.1.

The Nemenyi test results reveal that all pairwise comparisons between models show significant differences (p < 0.05). This finding corroborates the visual interpretation of Figure 4.4, confirming that the ranking distributions for all three models are statistically different from each other.

To further investigate potential differences in preferences between experts and novices in the domain of "Economy & Work", we employed an aligned rank transform (ART) ANOVA. This method was chosen for its ability to handle non-parametric data in factorial designs, making it particularly suitable for our ranking data. The ART ANOVA allows us to examine main effects and interactions while accounting for the within-subject nature of our design across different scenarios. Table 4.2 presents a comprehensive summary of the ART ANOVA results, including all main effects and interactions tested in our study.

The analysis revealed a highly significant main effect of model (F(2, 227) = 39.72, p < 0.001), confirming our earlier findings from the Friedman test. Interestingly, we found no significant main effect of Group (expert vs. novice) (F(1, 227) = 0, p = 1.000), indicating that

Factor/Interaction	Test Statistic	p-value	Significant
Model	F(2, 227) = 39.72	$10^{-15}$	Yes
Group (expert vs. novice)	F(1, 227) = 0	1.0	No
Scenario	F(9, 227) = 0	1.0	No
$Group \times Model$	F(2, 227) = 35.59	$10^{-14}$	Yes
Group $\times$ Scenario	F(9, 227) = 0	1.0	No
$Model \times Scenario$	F(18, 227) = 0.01	1.0	No

Table 4.2: ART ANOVA results showing significant effects for Model and Group  $\times$  Model interaction, with no significant effects for other factors or interactions. This indicates preference differences between experts and novices across summarization models, consistent across scenarios.

overall, experts and novices did not differ in their rankings across all models and scenarios. However, a significant Group  $\times$  Model interaction was observed (F(2,227)=35.59,p<0.001), suggesting that the ranking patterns for the three models differ between experts and novices.

To better understand this interaction, we conducted a detailed post-hoc analysis using Mann-Whitney U tests with Bonferroni correction for multiple comparisons. Table 4.3 presents the mean ranks and effect sizes for each model, comparing experts and novices.

While none of the individual comparisons between experts and novices for each model reached statistical significance after Bonferroni correction, the effect sizes suggest differences in how LLM Agent and UniChart are perceived by the two groups. Specifically, GPT-4V was consistently ranked highest by both experts and novices, with very little difference between the groups (effect size r = 0.0335). LLM Agent tended to be ranked lower by novices

Model	Mean Rank		Effect Size (r)	Adjusted p-value	Significant	
1,10,401	Expert	Novice	Effect Size (1)	riajustea p varae		
GPT-4V	1.45	1.50	0.0335	0.7289	No	
LLM Agent	1.98	2.30	0.1826	0.0720	No	
UniChart	2.57	2.20	0.1640	0.0870	No	

**Table 4.3:** Comparison of model rankings between experts and novices. Effect sizes (r) indicate the magnitude of difference between groups, with values closer to 0 suggesting smaller differences. Adjusted p-values are after Bonferroni correction.

compared to experts (mean ranks 2.30 vs. 1.98, effect size r = 0.1826). UniChart tended to be ranked lower by experts compared to novices (mean ranks 2.57 vs. 2.20, effect size r = 0.1640).

The lack of a significant scenario effect (F(9, 227) = 0, p = 1.000) in our ART ANOVA indicates that the performance of the models and the differences between expert and novice ratings were consistent across different types of charts. This suggests that the observed preferences are not dependent on specific chart types.

In conclusion, our analysis reveals differences in the perceived quality of chart summaries generated by different models, with GPT-4V consistently receiving the highest rankings across both experts and novices. The expertise level of the evaluators appears to influence the perception of LLM Agent and UniChart, but not GPT-4V, although these differences did not reach statistical significance in our post-hoc tests. Specifically, participants considered experts in the field of "Economy & Work" tended to rank LLM Agent higher, while novices tended to rank UniChart higher.

To understand the underlying reasons for these preferences, our next section examines the quality criteria, offering insights into each model's strengths and weaknesses for a more comprehensive evaluation of chart summarization performance.

### 4.5.2 Quality Criteria

Building upon our initial user preferences evaluation, this section presents a detailed analysis of model performance using the established quality criteria. As a brief reminder, the four key criteria are defined as follows: coherence measures the logical flow and organization of the summary; consistency evaluates the alignment between the summary and the original chart; fluency assesses the linguistic quality and readability of the summary; and relevance gauges how well the summary captures the essential information from the chart.

To compare the performance of GPT-4V, our LLM Agent, and UniChart across these criteria, we first conducted a descriptive analysis. Figure 4.5 presents the average scores and standard deviations for each model across all four criteria. The chart reveals that GPT-4V consistently achieved the highest average scores across all criteria, followed by our LLM Agent, with UniChart generally scoring the lowest. However, to determine the statistical significance of these differences, we employed a more rigorous statistical analysis.

Given our study design, which involved repeated measures (multiple evaluations by each participant) and ordinal data (Likert scale ratings), we chose the Friedman test as our initial statistical method. The Friedman test is particularly suitable for this scenario as it can

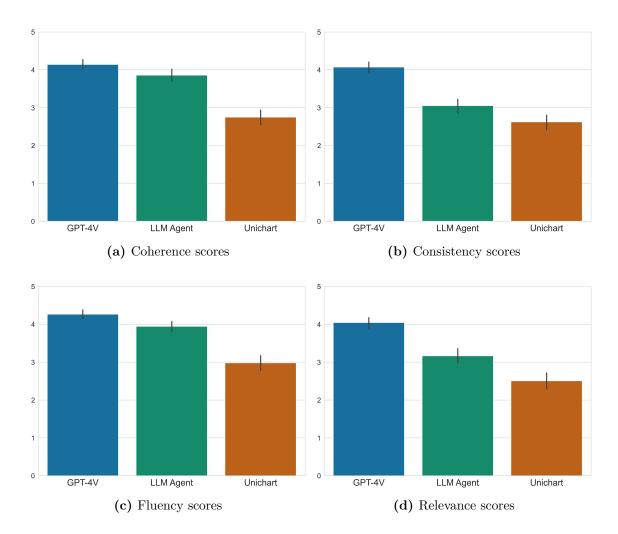


Figure 4.5: Comparison of Chart Summarization Model Performance Across Quality Criteria. This figure presents scores for GPT-4V, LLM Agent, and UniChart across four quality criteria: coherence, consistency, fluency, and relevance. GPT-4V consistently achieves the highest scores, followed by the LLM Agent, with UniChart generally scoring lowest. The performance gap is most pronounced in fluency and relevance. These results indicate GPT-4V's overall superiority, while highlighting the LLM Agent's competitive performance, particularly in linguistic aspects.

Criterion	Chi-square $(\chi^2)$	<i>p</i> -value	Significant
Coherence	17.3947	$10^{-4}$	Yes
Consistency	14.8462	$10^{-4}$	Yes
Fluency	24.3377	$10^{-6}$	Yes
Relevance	16.3000	$10^{-4}$	Yes

**Table 4.4:** Friedman Test Results for Differences in Quality Scores Between Models. This table presents the chi-square values, p-values, and significance for four key quality criteria: coherence, consistency, fluency, and relevance. All criteria show statistically significant differences (p < 0.001) between the models

detect differences between three or more related groups (in our case, GPT-4V, LLM Agent, and UniChart) without assuming normality in the data distribution. Table 4.4 presents the results of the Friedman test for each quality criterion.

The Friedman test results indicate statistically significant differences (p < 0.05) among the models for all four quality criteria. This confirms that the choice of model indeed impacts the quality of chart summaries across all evaluated aspects. To identify specific pairwise differences between models, we conducted a Nemenyi post-hoc test. We selected this test due to its conservative nature in controlling for family-wise error rates, which is necessary when performing multiple comparisons. Table 4.5 summarizes the results of the Nemenyi post-hoc test, including p-values, mean rank differences, and standardized differences for each model pair across all criteria.

GPT-4V emerges as the leader in chart summarization, consistently outperforming UniChart across all criteria with statistically significant differences (p < 0.05) and large

Criterion	Model Pair	<i>p</i> -value	Significant	Mean Rank Diff.	Std. Diff.
Coherence	GPT-4V vs. LLM Agent GPT-4V vs. UniChart LLM Agent vs. UniChart	0.646000 0.001000 0.007508	No Yes Yes	0.27500 1.22500 0.95000	0.86963 3.87379 3.00416
Consistency	GPT-4V vs. LLM Agent GPT-4V vs. UniChart LLM Agent vs. UniChart	0.083060 0.001000 0.220962	No Yes No	$0.67500 \\ 1.20000 \\ 0.52500$	2.13454 3.79473 1.66020
Fluency	GPT-4V vs. LLM Agent GPT-4V vs. UniChart LLM Agent vs. UniChart	0.099311 0.001000 0.015617	No Yes Yes	0.65000 1.52500 0.87500	2.05548 4.82247 2.76699
Relevance	GPT-4V vs. LLM Agent GPT-4V vs. UniChart LLM Agent vs. UniChart	0.019684 0.001000 0.416890	Yes Yes No	$0.85000 \\ 1.25000 \\ 0.40000$	2.68794 3.95285 1.26491

**Table 4.5:** Nemenyi Post-hoc Test Results for Pairwise Comparisons of Chart Summarization Models. The table shows statistical significance (p-value < 0.05), mean rank differences, and standardized differences across four quality criteria, quantifying performance gaps between GPT-4V, LLM Agent, and UniChart.

standardized mean rank differences. The gap between GPT-4V and UniChart is particularly pronounced, with standardized differences ranging from 3.79473 to 4.82247 across the four criteria. This substantial performance difference suggests that GPT-4V's multimodal capabilities provide a significant advantage in interpreting and summarizing chart data, allowing it to excel in both linguistic and data-centric aspects of the task.

Our LLM Agent demonstrates a varied performance profile, showing strengths in certain areas while facing challenges in others. In language-related criteria, specifically coherence and fluency, the LLM Agent significantly outperforms UniChart (p = 0.007508 and p = 0.015617, respectively), with large standardized differences of 3.00416 and 2.76699. This

strong performance in linguistic aspects indicates that our strategy of leveraging GPT-4's language capabilities in the agent's design has been effective. Moreover, the differences between GPT-4V and our LLM Agent are not statistically significant for these criteria (p = 0.646000 for coherence and p = 0.099311 for fluency), suggesting that our agent achieves comparable linguistic quality to GPT-4V in generating summaries.

However, the LLM Agent faces challenges in data-centric aspects of chart summarization. In terms of consistency and relevance, our agent shows no significant difference from UniChart (p=0.220962 and p=0.416890, respectively), with relatively small standardized differences of 1.66020 and 1.26491. This similarity in performance is not unexpected, as our agent utilizes UniChart's chart derendering and CVQA capabilities for information extraction (see section 3.2.3.2 on tools used by the action module). The contrast becomes more apparent when comparing the LLM Agent to GPT-4V in these areas. GPT-4V significantly outperforms our agent in relevance (p=0.019684, standardized difference = 2.68794) and shows a near-significant difference in consistency (p=0.083060, standardized difference = 2.13454). These results highlight that data-centric aspects are areas where GPT-4V's native multimodal capabilities provide a distinct advantage.

While the Friedman and Nemenyi tests provided valuable insights into the overall performance differences between models, they do not account for individual participant variability or the relative importance of each criterion in determining model rankings. To address these limitations and gain a more complete understanding of how experts and novices value chart summaries, we conducted an additional analysis using participant-level data.

This supplementary analysis employed two complementary statistical approaches: Spearman correlations with 95% confidence intervals and mixed-effects models. By using participant-level means for each model and criterion, we were able to preserve individual differences while reducing noise in the data. The Spearman correlations allow us to quantify the strength and direction of the relationship between each criterion and the overall model rankings, while the mixed-effects models account for both fixed effects (criteria) and random effects (participants). Table 4.6 presents the results of this analysis for both groups.

The results for the expert group demonstrate strong negative correlations between all criteria and model rankings, with values ranging from -0.618 to -0.724. The relatively small range of these correlation values suggests that experts value all four criteria almost equally in appreciating chart summaries. This finding is further supported by the narrow confidence intervals and highly significant p-values (all p < 0.001) obtained from the mixed-effects models.

Interestingly, while coherence shows the strongest overall correlation (-0.724), the mixed-effects model, which accounts for individual participant differences, reveals that relevance and consistency have the largest coefficients (-0.5210 and -0.5041, respectively).

Group	Criterion	Correlation	95% CI	Mixed-Model Coef.	p-value
Expert	Coherence	-0.724	[-0.890, -0.388] [-0.874, -0.325]	-0.402	1.01e-05 3.84e-14
	Consistency Fluency	-0.687 -0.618	[-0.842, -0.213]	-0.504 -0.458	3.63e-04
	Relevance	-0.665	[-0.864, -0.287]	-0.521	1.79e-26
Novice	Coherence	-0.116	[-0.848, 0.768]	-0.154	6.52 e-01
	Consistency	-0.721	[-0.967, 0.219]	-1.032	9.77e-02
	Fluency	-0.319	[-0.898, 0.665]	-0.319	5.91e-01
	Relevance	-0.588	[-0.948, 0.427]	-2.697	3.62e-12

Table 4.6: Spearman Correlation and Mixed-Effects Model Analysis Comparing Expert and Novice Evaluations of Summarization Models. The table presents correlation coefficients, 95% confidence intervals, mixed-model coefficients, and p-values across the four quality criteria, quantifying the relationship between criteria scores and overall model rankings for expert and novice evaluators.

This discrepancy suggests that when individual preferences are controlled for, data-centric criteria (relevance and consistency) emerge as more important in expert appreciations of chart summaries. However, the strong individual correlation for coherence indicates that some expert participants place high value on linguistic criteria.

The results for the novice group reveal a more complex picture. The correlations between criteria and rankings are generally weaker and not statistically significant, as evidenced by the wide confidence intervals that include zero. However, the mixed-effects model yields an intriguing result for relevance, showing a strong negative coefficient (-2.6970) with high statistical significance (p = 3.62e - 12).

This discrepancy between the correlation and mixed-effects model results for novices, particularly regarding relevance, warrants careful interpretation. The strong effect of relevance in the mixed-effects model suggests that when accounting for individual differences, relevance may play a crucial role in novice evaluations of chart summaries. This finding could explain our previous observation: novices do not differentiate in their preference between UniChart and our LLM Agent, as these models have similar relevance scores (see Table 4.1). In contrast, experts, who appear to value linguistic properties of chart summaries more highly, do distinguish between UniChart (with lower linguistic quality scores) and our LLM Agent (with higher linguistic scores). This pattern suggests that while novice users might primarily value data-centric criteria, particularly relevance, expert users consider both data-centric and linguistic aspects in their appreciations. However, it is important to note that the small sample size for the novice group (n=2) limits the reliability of these findings and necessitates further investigation with a larger cohort to confirm these patterns and the implied differences in evaluation priorities between novices and experts.

## 4.6 Discussion

This study aimed to evaluate the performance of our proposed LLM-based agent for chart summarization, compare it with existing baselines, and investigate the factors influencing user satisfaction across different expertise levels. Additionally, we sought to assess the effectiveness of our adapted quality criteria framework. The results of our user study and statistical analyses provide valuable insights into these research questions.

Addressing our first research question, the comparison between our LLM-based agent and the baselines (GPT-4V and UniChart) revealed a diverse performance profile. GPT-4V consistently outperformed both our agent and UniChart across all quality criteria, demonstrating the superiority of VLLMs in chart summarization tasks. However, our LLM-based agent showed promising results, particularly in linguistic aspects of summarization. The agent significantly outperformed UniChart in coherence and fluency, with no statistically significant difference from GPT-4V in these criteria. Importantly, our agent achieved these results without relying on a dataset of labeled chart summaries, which was a primary goal of this thesis. In terms of overall preference, our model successfully outperformed UniChart, validating our approach's effectiveness. These findings suggest that our strategy of leveraging GPT-4's language capabilities in the agent's design effectively enhanced the linguistic quality of the generated summaries while maintaining independence from costly labeled datasets.

Nevertheless, our LLM-based agent faced challenges in data-centric aspects of summarization, specifically consistency and relevance. In these criteria, our agent's performance was comparable to UniChart but significantly lower than GPT-4V. This limitation in data extraction and interpretation is not surprising, considering that our system relies on UniChart's CVQA and chart derendering components for these tasks. Consequently, our proposed method offers a tangible advantage: improvements in chart summarization performance can be achieved by training more accurate and powerful

CVQA and chart derendering models, without ever requiring costly labeled chart summaries.

Regarding our second research question on factors influencing user satisfaction, we observed notable differences between expert and novice evaluators. Experts demonstrated a balanced appreciation for all four quality criteria, with strong negative correlations between each criterion and model rankings. This suggests that experts value a holistic approach to chart summarization, considering both linguistic and data-centric aspects almost equally. Interestingly, when accounting for individual differences through mixed-effects modeling, relevance and consistency emerged as slightly more influential in expert evaluations. This finding indicates that while all expert participants shared a common emphasis on data-centric criteria, some individuals valued linguistic criteria particularly highly.

In contrast, novice evaluators showed a distinct preference pattern. While correlation analyses for novices were inconclusive due to small sample size, the mixed-effects model revealed a strong influence of relevance on their appreciations of the summaries. This finding suggests that novice users prioritize the inclusion of key information from the chart over linguistic sophistication. These differences between expert and novice preferences highlight the need for further investigation into how various user groups value different aspects of chart summaries, which could inform the development of more targeted and effective summarization systems.

The effectiveness of our adapted quality criteria framework is evident in its ability to

capture these fine differences in user preferences. By extending traditional text summarization criteria to the visual-textual domain of chart summarization, we were able to provide a more comprehensive and informative assessment of summary quality. The framework's ability to distinguish between linguistic and data-centric aspects of summaries proved particularly valuable in identifying the strengths and weaknesses of different models and understanding user preferences across different user groups.

Based on these findings, we propose several design recommendations for chart summarization systems. First, developers should focus on creating systems that can adapt to the needs of different user groups. The observed differences between experts and novices highlight a spectrum of user preferences that may vary across other groups as well. Further investigation using the quality criteria framework could reveal additional insights into these diverse needs. Second, with the advent of powerful text-based large language models like LLaMA 3, achieving GPT-4V level performance in chart summarization may be possible by focusing efforts on developing stronger CVQA and chart derendering tools. These components are likely easier to train than open-source VLLMs at the level of GPT-4, and improvements in these areas would directly enhance the performance of modular systems like our LLM-based agent. Finally, the integration of diverse quality criteria in the development and evaluation processes can lead to more robust and user-centric summarization systems, allowing for targeted improvements and more informative performance assessments.

# Chapter 5

## **Discussion**

## 5.1 Summary of Themes and Key Findings

This thesis investigates the application of LLM-based agents for chart summarization, focusing on enhancing accessibility for visually impaired individuals and addressing the scarcity of labeled datasets in this domain. The research yields several findings that contribute to the field of AI-assisted data interpretation and accessibility technologies.

A key theme that emerged is the efficacy of the LLM-based agent approach in chart summarization. By leveraging the general capabilities of LLMs in conjunction with domainspecific tools, this method demonstrates the potential to advance the field without relying on extensive labeled datasets, addressing a significant challenge in chart summarization research.

The study also introduces the application of quality criteria from text summarization to chart summarization. This approach provides a framework for evaluation, offering insights into the assessment of chart summary quality and user preferences. The criteria, which include coherence, consistency, fluency, and relevance, allow for a multifaceted evaluation of generated summaries. This framework proved valuable in differentiating model performance across linguistic and data-centric aspects, as well as in identifying preference patterns among different user groups.

The empirical evaluation compared the performance of the LLM-based agent to two baselines: GPT-4V and UniChart. GPT-4V consistently outperformed both the LLM-based agent and UniChart across all quality criteria. However, the LLM-based agent demonstrated competitive performance in linguistic aspects, specifically coherence and fluency, where it matched GPT-4V and outperformed UniChart. In data-centric aspects such as consistency and relevance, the LLM-based agent performed similarly to UniChart but fell short of GPT-4V's capabilities.

Analysis of user preferences revealed differences between expert and novice evaluators. Experts demonstrated a balanced appreciation for all quality criteria, valuing both linguistic and data-centric aspects of summaries. In contrast, novices placed higher emphasis on relevance. A trend emerged suggesting that experts tended to prefer the LLM agent's summaries, while novices leaned towards UniChart's outputs, although this difference did not reach statistical significance.

## 5.2 Implications and Practical Applications

The findings suggest a promising path forward for academia and open source to reach GPT-4 level performance, particularly as models like Meta's LLAMA-3-70B approach this benchmark. This strategy involves developing targeted, domain-specific models to enhance information extraction from charts, including improvements to chart derendering, CVQA, and new tools such as specialized OCR models.

The observed differences in preferences between expert and novice users reveal opportunities for customizable summary generation. In educational environments, systems could provide detailed explanations for novices while offering concise, insight-focused summaries for experts. This flexibility, achievable through prompt engineering in LLM-agents, could extend to various user categories, enabling personalized chart interpretation experiences.

Our quality criteria framework demonstrates value beyond evaluation, serving as a potential golden standard alongside user preference rankings in chart summarization research. These criteria not only enable precise performance comparison but also provide clear direction for future development efforts.

## 5.3 Limitations

A primary limitation of this study is the sample size and diversity of participants, particularly in the user evaluation phase. The focus on economic charts, while providing consistent context, may limit generalizability to other domains. Additionally, our evaluation primarily measured user preference rather than task-specific usefulness, limiting our understanding of practical effectiveness.

The current implementation faces technical constraints: increased processing time due to sequential GPT-4 API calls prevents real-time usage, and the static nature of summarization doesn't allow for interactive discovery. While our modular approach offers flexibility, it introduces system complexity and maintenance challenges.

### 5.4 Future Research Directions

To advance the field, we propose three key research directions:

First, enhancing evaluation through automated methods. Adapting frameworks like G-EVAL [68] could enable large-scale assessment of quality criteria, while employing LLMs for preference ranking could provide cost-effective evaluation. These methods would require validation against human judgments specifically for chart summarization.

Second, improving technical performance through enhanced tool utilization.

Implementing tree of thoughts [69] reasoning, integrating VLLMs as specialized tools, and

optimizing with smaller, distilled models like LLAMA-3 [70] could enhance both performance and efficiency.

Third, developing interactive summarization systems. Creating versions that allow follow-up questions could enhance user engagement and comprehension, particularly valuable for educational applications and professional analysis.

## 5.5 Conclusion

This thesis advances the field of AI-assisted data interpretation through its investigation of LLM-based agents for chart summarization. While our agent did not surpass GPT-4V, it demonstrated the viability of modular, tool-based approaches in addressing dataset scarcity. The developed quality criteria framework and insights into user preferences provide valuable foundations for future research.

The path forward involves improving data extraction capabilities, developing efficient LLM integration methods, and creating adaptive systems for different user needs. Most importantly, this work establishes a framework for developing accessible chart summarization systems that can serve diverse user groups while advancing the field through open-source and academic contributions.

- [1] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context. 2015. URL https://arxiv.org/abs/1405.0312.
- [2] Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. 2016. URL https://arxiv.org/abs/1505.04870.
- [3] Ahmed Masry, Parsa Kavehzadeh, Xuan Long Do, Enamul Hoque, and Shafiq Joty.
  UniChart: A Universal Vision-language Pretrained Model for Chart Comprehension
  and Reasoning. (arXiv:2305.14761), May 2023.
- [4] Shankar Kantharaj, Rixie Tiffany Ko Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq Joty. Chart-to-Text: A Large-Scale Benchmark for Chart Summarization. (arXiv:2203.06486), April 2022.

[5] Rin-ichiro Taniguchi, Masao Yokota, Eiji Kawaguchi, and Tuneo Tamati. Knowledge-based picture understanding of weather charts. *Pattern Recognition*, 17(1):109–123, January 1984. ISSN 0031-3203. doi: 10.1016/0031-3203(84)90038-4.

- [6] Weihua Huang, Chew Lim Tan, and Wee Kheng Leow. Associating text and graphics for scientific chart understanding. In *Eighth International Conference on Document* Analysis and Recognition (ICDAR'05), pages 580–584 Vol. 2, August 2005. doi: 10. 1109/ICDAR.2005.54.
- [7] Meredith Ringel Morris, Jazette Johnson, Cynthia L. Bennett, and Edward Cutrell. Rich Representations of Visual Content for Screen Reader Users. *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pages 1–11, April 2018. doi: 10.1145/3173574.3173633.
- [8] Rubén Alcaraz Martínez, Mireia Ribera Turró, and Toni Granollers Saltiveri. Methodology for heuristic evaluation of the accessibility of statistical charts for people with low vision and color vision deficiency. February 2021. ISSN 2693-5015. doi: 10.21203/rs.3.rs-156959/v1.
- [9] Anukriti Kumar, Tanuja Ganu, and Saikat Guha. ChartParser: Automatic Chart Parsing for Print-Impaired. (arXiv:2211.08863), November 2022. doi: 10.48550/arXiv. 2211.08863.
- [10] Prerna Mishra, Santosh Kumar, Mithilesh Kumar Chaube, and Urmila Shrawankar.

ChartVi: Charts summarizer for visually impaired. Journal of Computer Languages, 69:101107, April 2022. ISSN 2590-1184. doi: 10.1016/j.cola.2022.101107.

- [11] Fangyu Liu, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Yasemin Altun, Nigel Collier, and Julian Martin Eisenschlos. MatCha: Enhancing Visual Language Pretraining with Math Reasoning and Chart Derendering. (arXiv:2212.09662), December 2022.
- [12] R. Wang, C. Jung, and Y. Kim. Seeing Through Sounds: Mapping Auditory Dimensions to Data and Charts for People with Visual Impairments. Computer Graphics Forum, 41(3):71–83, 2022. ISSN 1467-8659. doi: 10.1111/cgf.14523.
- [13] Zeyuan Zhang. Tapsonic: One Dimensional Finger Mounted Multimodal Line Chart Reader. In *Proceedings of the 22nd International ACM SIGACCESS Conference on Computers and Accessibility*, ASSETS '20, pages 1–4, New York, NY, USA, October 2020. Association for Computing Machinery. ISBN 978-1-4503-7103-2. doi: 10.1145/3373625.3417075.
- [14] Nikolai A. Grabowski and Kenneth E. Barner. Data visualization methods for the blind using force feedback and sonification. In *Telemanipulator and Telepresence Technologies* V, volume 3524, pages 131–139. SPIE, December 1998. doi: 10.1117/12.333677.
- [15] Rayan Isran, Katayoun Sepehri, Karthikan Theivendran, and Ahmed Anwar. Towards

  More Effective Data Visualization Methods Using Haptics. In 2021 IEEE World Haptics

Conference (WHC), pages 590–590, Montreal, QC, Canada, July 2021. IEEE. ISBN 978-1-66541-871-3. doi: 10.1109/WHC49131.2021.9517255.

- [16] Jiho Kim, Arjun Srinivasan, Nam Wook Kim, and Yea-Seul Kim. Exploring Chart Question Answering for Blind and Low Vision Users. In *Proceedings of the 2023 CHI* Conference on Human Factors in Computing Systems, CHI '23, pages 1–15, New York, NY, USA, April 2023. Association for Computing Machinery. ISBN 978-1-4503-9421-5. doi: 10.1145/3544548.3581532.
- [17] Filip Bajić and Josip Job. Review of chart image detection and classification. International Journal on Document Analysis and Recognition (IJDAR), January 2023. ISSN 1433-2825. doi: 10.1007/s10032-022-00424-5.
- [18] Himanshu Sharma and Devanand Padha. A comprehensive survey on image captioning: From handcrafted to deep learning-based techniques, a taxonomy and open research issues. *Artificial Intelligence Review*, 56(11):13619–13661, November 2023. ISSN 1573-7462. doi: 10.1007/s10462-023-10488-2.
- [19] Girish Kulkarni, Visruth Premraj, Vicente Ordonez, Sagnik Dhar, Siming Li, Yejin Choi, Alexander C. Berg, and Tamara L. Berg. BabyTalk: Understanding and Generating Simple Image Descriptions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2891–2903, December 2013. ISSN 1939-3539. doi: 10.1109/TPAMI. 2012.162.

[20] Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2Text: Describing Images Using 1 Million Captioned Photographs. In Advances in Neural Information Processing Systems, volume 24. Curran Associates, Inc., 2011.

- [21] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation Applied to Handwritten Zip Code Recognition. *Neural Computation*, 1(4):541–551, December 1989. ISSN 0899-7667. doi: 10.1162/neco.1989. 1.4.541.
- [22] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, October 1986. ISSN 0028-0836, 1476-4687. doi: 10.1038/323533a0.
- [23] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and Tell: A Neural Image Caption Generator. (arXiv:1411.4555), April 2015. doi: 10.48550/arXiv. 1411.4555.
- [24] Hrituraj Singh and Sumit Shekhar. STL-CQA: Structure-based Transformers with Localization and Encoding for Chart Question Answering. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3275–3284, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020. emnlp-main.264.

[25] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, and Yoshua Bengio. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. (arXiv:1502.03044), April 2016. doi: 10.48550/arXiv.1502.03044.

- [26] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. ChartQA:
  A Benchmark for Question Answering about Charts with Visual and Logical Reasoning.
  (arXiv:2203.10244), March 2022.
- [27] Fangyu Liu, Julian Martin Eisenschlos, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Wenhu Chen, Nigel Collier, and Yasemin Altun. DePlot: One-shot visual language reasoning by plot-to-table translation.

  (arXiv:2212.10505), December 2022.
- [28] Teodor Fredriksson, David Issa Mattos, Jan Bosch, and Helena Holmström Olsson. Data labeling: An empirical investigation into industrial challenges and mitigation strategies. In Maurizio Morisio, Marco Torchiano, and Andreas Jedlitschka, editors, Product-Focused Software Process Improvement, pages 202–216, Cham, 2020. Springer International Publishing. ISBN 978-3-030-64148-1.
- [29] Saad Obaid ul Islam, Iza Škrjanec, Ondřej Dušek, and Vera Demberg. Tackling Hallucinations in Neural Chart Summarization. (arXiv:2308.00399), August 2023. doi: 10.48550/arXiv.2308.00399.

[30] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. (arXiv:2203.02155), March 2022. doi: 10.48550/arXiv.2203.02155.

- [31] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling Instruction-Finetuned Language Models. (arXiv:2210.11416), December 2022.
- [32] Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. DVQA: Understanding Data Visualizations via Question Answering. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 5648–5656, 2018.
- [33] Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Akos Kadar, Adam

Trischler, and Yoshua Bengio. FigureQA: An Annotated Figure Dataset for Visual Reasoning. (arXiv:1710.07300), February 2018.

- [34] Ritwick Chaudhry, Sumit Shekhar, Utkarsh Gupta, Pranav Maneriker, Prann Bansal, and Ajay Joshi. LEAF-QA: Locate, Encode & Attend for Figure Question Answering. (arXiv:1907.12861), July 2019. doi: 10.48550/arXiv.1907.12861.
- [35] Nitesh Methani, Pritha Ganguly, Mitesh M. Khapra, and Pratyush Kumar. PlotQA: Reasoning over Scientific Plots. (arXiv:1909.00997), February 2020.
- [36] Shankar Kantharaj, Xuan Long Do, Rixie Tiffany Ko Leong, Jia Qing Tan, Enamul Hoque, and Shafiq Joty. OpenCQA: Open-ended Question Answering with Charts. (arXiv:2210.06628), October 2022. doi: 10.48550/arXiv.2210.06628.
- [37] Kanit Wongsuphasawat, Zening Qu, Dominik Moritz, Riley Chang, Felix Ouk, Anushka Anand, Jock Mackinlay, Bill Howe, and Jeffrey Heer. Voyager 2: Augmenting Visual Analysis with Partial View Specifications. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, CHI '17, pages 2648–2659, New York, NY, USA, May 2017. Association for Computing Machinery. ISBN 978-1-4503-4655-9. doi: 10.1145/3025453.3025768.
- [38] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu,

Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. The Rise and Potential of Large Language Model Based Agents: A Survey. (arXiv:2309.07864), September 2023.

- [39] OpenAI and Others. GPT-4 Technical Report. (arXiv:2303.08774), March 2024. doi: 10.48550/arXiv.2303.08774. Full author list includes 280 contributors.
- [40] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large Language Models are Zero-Shot Reasoners. (arXiv:2205.11916), January 2023. doi: 10.48550/arXiv.2205.11916.
- [41] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. (arXiv:2005.14165), July 2020. doi: 10.48550/arXiv. 2005.14165.
- [42] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. (arXiv:2201.11903), January 2023.

[43] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-Consistency Improves Chain of Thought Reasoning in Language Models. (arXiv:2203.11171), March 2023. doi: 10.48550/arXiv. 2203.11171.

- [44] Albert Lu, Hongxin Zhang, Yanzhe Zhang, Xuezhi Wang, and Diyi Yang. Bounding the Capabilities of Large Language Models in Open Text Generation with Prompt Constraints. (arXiv:2302.09185), February 2023. doi: 10.48550/arXiv.2302.09185.
- [45] Nick McKenna, Tianyi Li, Liang Cheng, Mohammad Javad Hosseini, Mark Johnson, and Mark Steedman. Sources of Hallucination by Large Language Models on Inference Tasks. (arXiv:2305.14552), October 2023. doi: 10.48550/arXiv.2305.14552.
- [46] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language Agents with Verbal Reinforcement Learning. (arXiv:2303.11366), October 2023. doi: 10.48550/arXiv.2303.11366.
- [47] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. (arXiv:2010.11929), June 2021.
- [48] OpenAI. Gpt-4 vision system card. September 2023. URL https://cdn.openai.com/papers/GPTV\_System\_Card.pdf.

[49] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. (arXiv:2301.12597), January 2023.

- [50] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. ReAct: Synergizing Reasoning and Acting in Language Models. (arXiv:2210.03629), March 2023.
- [51] Jie Huang and Kevin Chen-Chuan Chang. Towards Reasoning in Large Language Models: A Survey. (arXiv:2212.10403), May 2023. doi: 10.48550/arXiv.2212.10403.
- [52] Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language Models can Solve Computer Tasks. (arXiv:2303.17491), November 2023. doi: 10.48550/arXiv.2303.17491.
- [53] Yunjia Zhang, Jordan Henkel, Avrilia Floratou, Joyce Cahoon, Shaleen Deep, and Jignesh M. Patel. ReAcTable: Enhancing ReAct for Table Question Answering. (arXiv:2310.00815), October 2023. doi: 10.48550/arXiv.2310.00815.
- [54] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. ToolLLM: Facilitating Large Language Models to Master 16000+ Real-world APIs. (arXiv:2307.16789), July 2023.
- [55] Shikun Liu, Linxi Fan, Edward Johns, Zhiding Yu, Chaowei Xiao, and Anima

Anandkumar. Prismer: A Vision-Language Model with An Ensemble of Experts. (arXiv:2303.02506), March 2023.

- [56] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual ChatGPT: Talking, Drawing and Editing with Visual Foundation Models. (arXiv:2303.04671), March 2023.
- [57] Wojciech Kryscinski, Nitish Shirish Keskar, Bryan McCann, Caiming Xiong, and Richard Socher. Neural Text Summarization: A Critical Evaluation. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 540–551, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1051.
- [58] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. BERTScore: Evaluating Text Generation with BERT. (arXiv:1904.09675), February 2020. doi: 10.48550/arXiv.1904.09675.
- [59] Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. SummEval: Re-evaluating Summarization Evaluation. Transactions of the Association for Computational Linguistics, 9:391–409, April 2021. ISSN 2307-387X. doi: 10.1162/tacl\_a\_00373.

[60] Divakar Yadav, Jalpa Desai, and Arun Kumar Yadav. Automatic Text Summarization Methods: A Comprehensive Review. (arXiv:2204.01849), March 2022. doi: 10.48550/ arXiv.2204.01849.

- [61] OpenAI. OpenAI Models Documentation. URL https://platform.openai.com/ docs/models.
- [62] Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. Berkeley function calling leaderboard. 2024.
- [63] Walter Kintsch and Teun A. van Dijk. Toward a model of text comprehension and production. Psychological Review, 85(5):363–394, 1978. ISSN 1939-1471. doi: 10.1037/ 0033-295X.85.5.363.
- [64] George A. Miller. The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological Review*, 63(2):81–97, 1956. ISSN 1939-1471. doi: 10.1037/h0043158.
- [65] Haley MacLeod, Cynthia L. Bennett, Meredith Ringel Morris, and Edward Cutrell. Understanding Blind People's Experiences with Computer-Generated Captions of Social Media Images. In Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems, CHI '17, pages 5988–5999, New York, NY, USA, May 2017. Association for Computing Machinery. ISBN 978-1-4503-4655-9. doi: 10.1145/3025453. 3025814.

[66] Ori Ernst, Ori Shapira, Ido Dagan, and Ran Levy. Re-Examining Summarization Evaluation across Multiple Quality Criteria. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, Findings of the Association for Computational Linguistics: EMNLP 2023, pages 13829–13838, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.924.

- [67] Dan Gillick and Yang Liu. Non-Expert Evaluation of Summarization Systems is Risky. In Chris Callison-Burch and Mark Dredze, editors, *Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon's Mechanical Turk*, pages 148–151, Los Angeles, June 2010. Association for Computational Linguistics.
- [68] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-Eval: NLG Evaluation using GPT-4 with Better Human Alignment. (arXiv:2303.16634), May 2023. doi: 10.48550/arXiv.2303.16634.
- [69] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of Thoughts: Deliberate Problem Solving with Large Language Models. (arXiv:2305.10601), May 2023.
- [70] Abhimanyu Dubey and Others. The llama 3 herd of models. 2024. URL https://arxiv.org/abs/2407.21783. Full author list includes 280 contributors.

# Appendix A

## Expertise Questionnaire

- 1. Which of the following best describes the relationship between inflation and unemployment rates in the short run? (Medium)
  - a) No relationship
  - b) Direct, positive relationship
  - c) Inverse relationship, as described by the Phillips Curve
  - d) Always move in the same direction
- 2. How does the Consumer Price Index (CPI) measure inflation? (Easy)
  - a) Tracking price changes in a fixed basket of goods/services
  - b) Measuring quarterly changes in GDP
  - c) Calculating average price of a changing basket of goods/services
  - d) Surveying consumers on perceived living costs

3.	Which of the following is a potential cause of cost-push inflation? (Medium)
	a) Increasing consumer demand
	b) Rising wages due to labor shortages
	c) Expansionary monetary policy
	d) Increasing global competition
4.	Which of the following is NOT a component of GDP? (Easy)
	a) Consumer spending
	b) Government spending
	c) Stock market performance
	d) Net exports
<b>5</b> .	Central banks can indirectly influence which of the following? (Hard)
	a) Government tax rates
	b) Bank lending rates
	c) Corporate hiring decisions
	d) Stock market valuations
6.	Which of the following is an example of a progressive tax system? (Medium)
	a) A flat income tax rate for all taxpayers
	b) Higher income tax rates for lower-income earners

c) Lower sales tax rates on luxury goods

- d) Higher income tax rates for higher-income earners
- 7. The unemployment rate is calculated as: (Easy)
  - a) # unemployed / total population
  - b) # unemployed / labor force
  - c) # employed / labor force
  - d) # unemployed / # employed
- 8. Which of the following would be considered a discouraged worker? (Hard)
  - a) A person who is actively seeking work but unable to find a job that matches their qualifications
  - b) A person who wants a job but has given up looking due to lack of opportunities
  - c) A person who is working part-time but wants to work full-time for better pay and benefits
  - d) A person who is unemployed and not actively seeking work due to family responsibilities
- 9. A bear market refers to: (Easy)
  - a) A period of rising stock prices
  - b) A period of falling stock prices
  - c) High volatility in the housing market

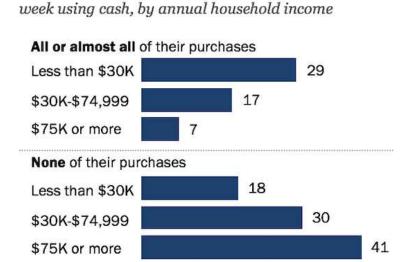
- d) Rapid growth in the money supply
- 10. Which of the following best describes the primary goal of monetary policy? (Medium)
  - a) Achieving price stability and sustainable economic growth
  - b) Balancing the government budget
  - c) Increasing government spending on infrastructure
  - d) Encouraging foreign investment in domestic markets
- 11. Which of the following is generally considered a leading indicator for the housing market?

  (Medium)
  - a) New home sales
  - b) Existing home sales
  - c) Housing starts
  - d) Home price appreciation
- 12. Which of the following is most likely to contribute to increasing wealth inequality? (Hard)
  - a) Regressive taxation policies that benefit high-income earners
  - b) Increasing prevalence of part-time and contract work arrangements
  - c) Stock market gains accruing primarily to the wealthy
  - d) Uneven access to healthcare and health insurance based on socioeconomic status
- 13. The labor force participation rate measures: (Medium)
  - a) The percentage of the total population that is employed

- b) The percentage of the working-age population that is either employed or actively seeking work
- c) The percentage of the working-age population that is employed
- d) The percentage of the total population that is unemployed
- 14. Which of the following best describes the impact of rising interest rates on the housing market? (Hard)
  - a) Encourages more people to buy homes before interest rates increase further
  - b) Decreases affordability and demand for homes
  - c) Leads to a higher supply of homes as builders rush to complete projects
  - d) Increases the number of adjustable-rate mortgages as buyers seek lower initial payments
- **15.** Which of the following is an example of expansionary fiscal policy? (Medium)
  - a) Raising taxes
  - b) Cutting government spending
  - c) Increasing government spending on infrastructure projects
  - d) The central bank selling government bonds

# Appendix B

## Summary Generation Example



% of U.S. adults who say they make \_\_\_\_ in a typical

Note: Respondents who did not give an answer or gave other responses are not shown.

Source: Survey conducted Sept. 24-Oct. 7, 2018.

Figure B.1: Percentage of U.S. adults' cash purchases by income level

## Chart Analysis

This appendix presents an example of our LLM-based agent generating a summary for the chart on Figure B.1.

```
Action: Calling 'pandas-code-interpreter'
Arguments:
print (df)
Observation:
                             Characteristic
                                               Value
   All or almost all of their purchases
                                                0.29
1
                             Less than $30K
                                                0.29
2
                              $30K-$74,999
                                                0.17
3
                              $75K or more
                                                0.07
4
                  None of their purchases
                                                0.18
5
                            Less than $30K
                                                0.18
6
                               $30K-$74.999
                                                0.30
7
                               $75K or more
                                                0.41
```

### Thought

The dataframe head provides a glimpse into the data that the chart is likely based on. It appears to be related to some form of purchasing behavior or financial data, given the references to "All or almost all of their purchases" and income brackets such as "Less than \$30K," "\$30K-\$74,999," and "\$75K or more." The 'Value' column seems to represent a percentage or a proportion, as suggested by the decimal format, which is typical for such data.

To better understand the chart, I need to establish the following: 1. The overall topic or subject of the chart. 2. The specific entities or categories being compared or tracked. 3. The range and distribution of values to understand the scale and significance of differences. 4. Any trends or patterns that might be evident from the data. 5. The context or demographic information that might explain the purchasing behavior or financial data.

#### Action

Action: Calling 'unichart-opencqa'

**Arguments:** question = "What is the chart about?"

**Observation:** Around three-in-ten Americans (29%) say they make less than \$30,000 a year ago, making those earning less than \$75,000 or more a year ago. Only 6% of Americans say they have purchases less than \$30,000 a year ago, making those earning more than \$75,000 or more a year ago.

## Thought

While this gives us some initial context about income levels and purchasing behavior, we need to understand the specific nature of these purchases and when this data was collected. This temporal and contextual information will be crucial for properly interpreting the patterns we're seeing.

#### Action

Action: Calling 'unichart-opencqa'

**Arguments:** question = "What is the specific context of these purchases and when was this data collected?"

**Observation:** This data shows the percentage of U.S. adults who make purchases using cash in a typical week, based on a survey conducted from September 24 to October 7, 2018. The data excludes respondents who did not provide an answer or gave other responses.

## Thought

The responses from the 'unichart-opencqa' functions have provided valuable context. The chart examines cash purchasing behavior among Americans, based on a 2018 survey. The specific focus on cash transactions adds an important dimension to our understanding. To further analyze the relationships shown in the chart, I need to: 1. Confirm the categories or groups being compared in the chart. 2. Understand the specific values associated with each category. 3. Determine the relationship between income levels and cash purchasing behavior. 4. Identify any trends or patterns in the data, such as whether higher income correlates with certain cash payment patterns.

#### Action

Action: Calling 'unichart-closedga'

## Arguments & Observations:

Q: "What are the categories being compared in the chart?"

 $\rightarrow$  [All or almost all of their purchases, Less than \$30K]

Q: "What is the value associated with the category 'Less than \$30K'?"

 $\rightarrow 29$ 

Q: "What is the value associated with the category '\$75K or more'?"

 $\rightarrow 41$ 

## Thought

The responses from the 'unichart-closedqa' function have provided some specific details. The categories being compared in the chart include "All or almost all of their purchases" and income levels such as "Less than \$30K." The value associated with the category "Less than \$30K" is 29, and the value for the category "\$75K or more" is 41. These values represent percentages of cash purchases, given the context provided earlier.

With this information, it is becoming clearer that the chart is comparing the percentage of Americans in different income brackets and their cash payment behaviors. The data suggests a possible correlation between income levels and cash usage patterns, with higher income brackets showing different tendencies.

#### Action

**Action:** Calling 'pandas-code-interpreter'

### **Arguments:**

```
print(df['Value'].mean())
print(df['Value'].std())
```

### Observation:

0.23624999999999999

0.10582161539928542

## Thought

**Final Anwer:** The chart appears to be comparing the cash purchasing behavior of Americans across different income levels, based on a 2018 survey. The dataframe reveals four categories of purchasing behavior: "All or almost all of their purchases" and "None of their purchases," with corresponding income brackets "Less than \$30K," "\$30K-\$74,999," and "\$75K or more."

The mean value across these categories is approximately 23.62%, with a standard deviation of about 10.58%. This suggests that there is some variation in the cash purchasing behavior across different income levels, but not an extremely wide dispersion.

The specific values for each category are as follows:

- "All or almost all of their purchases" using cash corresponds to 29% for "Less than \$30K," 17% for "\$30K-\$74,999," and 7% for "\$75K or more."
- "None of their purchases" using cash corresponds to 18% for "Less than \$30K," 30% for "\$30K-\$74,999," and 41% for "\$75K or more."

From this data, we can infer that as income increases, the percentage of people who make "All or almost all of their purchases" using cash decreases, while the percentage of people who make "None of their purchases" using cash increases. This could suggest that higher-income individuals are less likely to use cash as their primary payment method.

## Generated Summary

Based on the extracted information and analysis above, the system generated the following reasoning and summary. While both are shown here for completeness, only the final summary would be provided to the user:

## Summary

Reasoning: The chart data indicates a relationship between Americans' income levels and their cash purchasing behavior, based on a 2018 survey. The categories include 'All or almost all of their purchases' and 'None of their purchases,' across three income brackets: 'Less than \$30K,' '\$30K-\$74,999,' and '\$75K or more.' The mean value of 23.62% and standard deviation of 10.58% suggest moderate variation across categories. Higher income correlates with a lower percentage of 'All or almost all of their purchases' using cash (7% for '\$75K or more') and a higher percentage of 'None of their purchases' using cash (41% for '\$75K or more'). This trend inverses for lower income levels, indicating that income influences cash payment preferences.

**Summary:** The chart illustrates a correlation between income levels and cash payment behavior in America as of 2018. Higher-income earners are less likely to make 'All or almost all of their purchases' using cash compared to lower-income earners, who report a higher percentage of such purchasing behavior. Conversely, higher-income individuals report a greater percentage of 'None of their purchases' using cash, suggesting income-based disparities in payment method preferences.