# ChartCheck: Explainable Fact-Checking over Real-World Chart Images

**Mubashara Akhtar[1], Nikesh Subedi[2], Vivek Gupta[3], Sahar Tahmasebi[4],**
**Oana Cocarascu[1] and Elena Simperl[1]**

[1]King's College London [2]University of Utah [3]University of Pennsylvania

[4]TIB – Leibniz Information Centre for Science and Technology

`mubashara.akhtar@kcl.ac.uk`

## Abstract

Whilst fact verification has attracted substantial interest in the natural language processing community, verifying misinforming statements against data visualizations such as charts has so far been overlooked. Charts are commonly used in the real-world to summarize and communicate key information, but they can also be easily misused to spread misinformation and promote certain agendas. In this paper, we introduce ChartCheck, a novel, large-scale dataset for explainable fact-checking against real-world charts, consisting of 1.7k charts and 10.5k human-written claims and explanations. We systematically evaluate ChartCheck using vision-language and chart-to-table models, and propose a baseline to the community. Finally, we study chart reasoning types and visual attributes that pose a challenge to these models.[1]

## 1 Introduction

Data visualizations (e.g. bar charts, pie charts, line graphs) are common in real-world data sources and are frequently used in scientific documents, textbooks, news articles, and social media to summarize key information in a visual form, information which is often not fully repeated in the associated text (Carberry et al., 2006). However, data visualizations are also commonly misused to spread misinformation,[2] being included in political and business advertisements, or spread on social media, with the aim to convince the audience towards certain agendas (Lo et al., 2022; Lisnic et al., 2023). For example, during the COVID pandemic, the term *"counter-visualizations"* was coined for charts circulating on social media platforms to counter COVID-related health measures (Lee et al., 2021) (see Fig. 2). Charts were also (mis-)used by Brexit campaigners during the referendum.[3]



**Evidence**

**Chart:**

2006 Mexican Presidential Election Vote Count Progression

**Caption:** 2006 Mexican Presidential election vote count progression. Percentage of poll stations counted vs. percentage of candidate votes.

**Claim:** The percentage of votes for Loprez Obrador decreased over time as more poll stations were counted.
**Explanation:** The chart shows a generally downward trend in the percentage of votes for Obrador as more poll stations were counted. As more votes were counted, the other candidate gained a larger percentage of the overall vote.
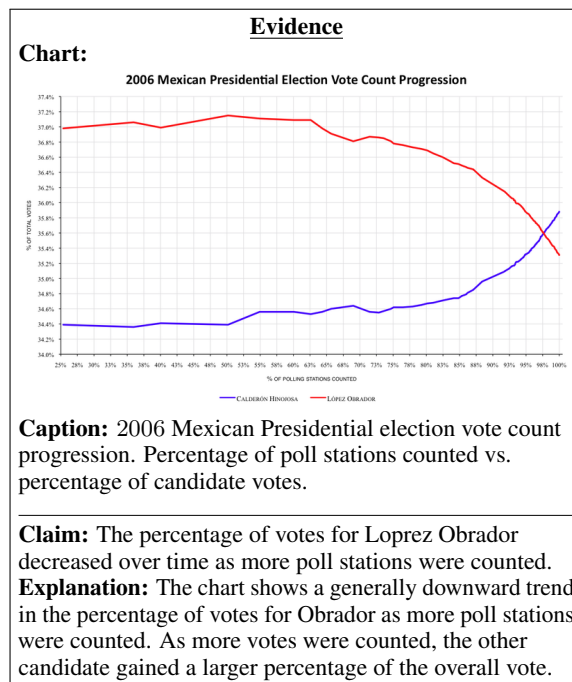
Figure 1: A dataset instance labelled as *supported*. Example includes the claim, the chart, its caption, and an explanation justifying why the claim is supported.

While previous research has focused on misleading chart design (i.e. truncated axis) (Lo et al., 2024; Hemsley and Snyder, 2018), a less studied but prevalent issue is misleading statements that exploit information fallacies during chart interpretation. Verifying these statements requires extracting information from charts that use tightly integrated text and visual elements to represent information. These visual elements consist of various lines and shapes, and have different colors, scales, angles, and orientations. Moreover, various reasoning skills, including numeracy and language understanding, are required to successfully verify statements against charts.

In this paper, we introduce *ChartCheck*, a novel dataset for explainable fact-checking against data visualizations. ChartCheck is the first large-scale dataset with 1.7*k* real-world charts extracted from

---

[1]Code and data available at `https://github.com/mubasharaak/ChartCheck`
[2]See Fig. 2 and appendix for examples.
[3]See Fig. 10 and 11 in the appendix.

the web and $10.5k$ human-written claims and explanations (see Fig. 1 for an example). Different to previous datasets, ChartCheck provides explanations as justifications for claim verification. Explaining fact-checking decisions is an important aspect as a major goal of fact-checking is to convince readers, and debunking misinformation by only labelling it as false is often not sufficient (Guo et al., 2022). To collect the data, we apply a four-step crowdsourcing pipeline to filter out noisy images, write claims and explanations, evaluate the data, and apply a final verification check through expert annotators.

We evaluate state-of-the-art (SOTA) models in a finetuning, few- and zero-shot setting. Hereby, we consider two baseline architectures: (1) vision-language models (VLMs) and (2) chart-to-table architectures with separate models for each subproblem: chart translation (chart-to-table), claim verification, and explanation generation. Our best-performing chart-to-table baseline reaches 73.8 accuracy and lags far behind human performance.

We annotate a subset of ChartCheck with $(i)$ reasoning types humans apply for chart understanding (Amar et al., 2005) and $(ii)$ visual attributes common in charts. We evaluate which reasoning types and chart attributes pose a challenge for SOTA models, as well as common failure types related to explanation generation. We evaluate model-generate explanations for factuality, semantic coverage, coherence, readability, and redundancy. Overall, our results show that models generate very fluent, readable and coherent text, but show major limitations w.r.t. understanding and reasoning over charts to generate correct explanations. Hence, ChartCheck is a challenging problem and will stimulate progress on fact-checking against charts.

To summarise, our **contributions** are as follows:
*1)* We propose ChartCheck, the first chart dataset for explainable fact-verification with $1.7k$ real-world charts and $10.5k$ human-written claims and explanations.
*2)* We evaluate chart-to-table models and VLMs in a finetuning and few-/zero-shot setting.
*3)* We study chart reasoning types and visual attributes that pose a challenge to SOTA models, as well as failure types related to explanations.

## 2   The Need for Chart Fact-checking

**Charts in the real-world.** Data visualizations are commonly consumed on a daily basis, e.g. in news articles, text books, scientific papers, and on the In-
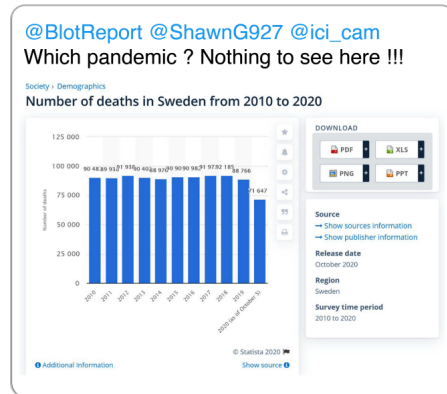


Figure 2: Twitter post related to COVID misinformation claiming that number of deaths in Sweden did not increase through the pandemic although data for the year 2020 is not complete (Lisnic et al., 2023).

ternet (Lo et al., 2022). Charts, diagrams, plots and infographics are useful tools to summarize and communicate key information in a visual form, which is often not fully repeated in the associated text (Carberry et al., 2006). Moreover, data visualizations are popular data sources for studying data-centric problems in different domains such as finance, science, health, climate-change. For example, during the pandemic, charts were widely used to guide policymakers in deciding health policies and communicate COVID information to the general public (Johns Hopkins University's coronavirus dashboard is a popular example).[4]

**Misinformed by charts.** Despite their useful applications, data visualizations are also commonly misused to spread misinformation. Previous research on chart-related misinformation focused on misinformation caused by misleading chart design, in particular on visual manipulation techniques, such as truncated or double axes, missing legends, linear scale for exponential data, confusing colors, or violation of guidelines and best practises related to visualization design (Lo et al., 2022). However, a misinformation type related to data visualizations that is less studied but prevalent in the real-world, are information fallacies during chart interpretation. Analyzing Twitter posts that contain data visualizations, Lisnic et al. (2023) found that the most common way people mislead with charts is through misleading arguments that exploit information fallacies as opposed to misleading chart design. They identify common components of misleading chart interpretation such as cherry-picking of the data points and time frames, setting of arbitrary data

---

[4] https://coronavirus.jhu.edu/map.html

thresholds, and causal errors, among others.

**Reasoning over charts.** Misinformation related to chart interpretations cannot be identified through visual manipulation detectors but requires multiple sub-steps: $(i)$ understanding the chart content given its context (e.g. caption, surrounding information); $(ii)$ understanding the accompanying claim; $(iii)$ identifying and performing necessary reasoning steps (e.g. arithmetic); $(iv)$ deciding and explaining why the claim is correct or not. There are certain challenges unique to reasoning and verification over data visualizations. To verify a claim, information needs to be extract from chart images that convey meaning through various types of visual elements such as lines, shapes, different colors, scales, angles, orientations, etc. (Lee et al., 2022; Liu et al., 2022b). Moreover, various skills, including numeracy, causal reasoning, and understanding of semantics, need to be successfully applied to verify the given claim. While there are many prior works on visual-language misinformation, these have mostly focused on datasets with natural images where the language and visual information is not strongly integrated, e.g. manipulation or out-of-context detection (Akhtar et al., 2023b).

## 3 ChartCheck: Dataset Creation

Fig. 3 provides an overview of the dataset creation pipeline. Starting with chart extraction from the web (step 1), we applied three crowdsourcing tasks (step $2-4$), before post-processing and validating the collected data (step 5).

### 3.1 Data Collection (Step 1)

To collect a diverse set of chart images, we used Wikimedia Commons.[5] We extracted horizontal/vertical bar charts, line/area charts, pie/donuts charts, and scatterplots, if the chart description was in English, resulting in $2,498$ charts.[6]

### 3.2 Crowdsourcing (Steps $2-4$)

We used Mechanical Turk for crowdsourcing.[7]

**Chart filtering (Step** 2**).** First, we filter out noisy chart images for the subsequent tasks. We asked annotators to label if charts were $(i)$ non-English, $(ii)$ not readable, $(iii)$ not understandable (e.g., because crucial information such as the legend was

---

5Details on worker training/recruitment in Appendix C.2.

[5]https://commons.wikimedia.org/wiki/Commons:Welcome

[6]See Appendix C.1 for further details.

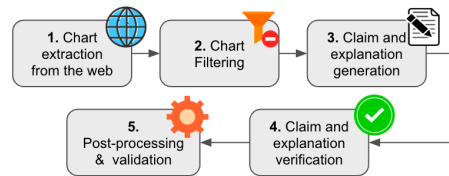[7]Details on worker training/recruitment in Appendix C.2.



Figure 3: Overview of the ChartCheck dataset pipeline.

missing). We implemented automated checks required to submit a task, e.g. if all labels were set and annotators spent a minimum time span on each chart. Moreover, we included in each taskset of seven chart images two golden samples pre-labeled by the authors. Annotators could only submit a taskset if they labelled the gold samples correctly. Finally, we filtered the collected data using majority voting, resulting in $1,684$ charts.

**Claim and explanation generation (Step** 3**).** Claims and explanations were written by a different group of annotators. For each chart, we asked annotators to write one claim supported by the chart and one refuted. The workers also wrote a text explaining why the claim is supported/refuted and outlining the reasoning steps. To improve the submission quality and guarantee non-superficial claims and explanations, we implemented automated checks and a manual verification step. We automatically evaluated the submissions and rejected those that violated the following conditions: $(i)$ claims had a length between $7-30$ words and explanations between $10-100$ words; $(ii)$ no duplicate claims or explanations were submitted; $(iii)$ claims/explanations did not contain uncertainty terms (e.g. *"perhaps"*, *"maybe"*, *"occasionally"*); $(iv)$ refuting claims were not simple negations of supporting claim; $(v)$ minimum time spent on each chart was $5$ seconds. Moreover, we instructed workers to refrain from writing simple claims and gave corresponding examples. Finally, we randomly sampled and manually evaluated the quality of one claims-explanation pair for each submitted task before accepting the annotation work. In total, we collected $9,300$ claims and explanations. On average, claims and explanations have a length of $12.5$ and $23.4$ words, respectively.

**Validation of claims and explanations (Step** 4**).** Next, we asked another group of workers to verify the claims and explanations. If the assigned claim label did not match its content, workers corrected them. Similarly, workers corrected explanations if they did not explain correctly why a claim was supported/refuted by the chart and its

| Set | #Images | #Claims | Support | Refute |
|---|---|---|---|---|
| **Train** | 1,532 | 7,607 | 3,871 | 3,736 |
| **Valid** | 672 | 953 | 494 | 459 |
| **Test** $t1$ | 669 | 939 | 487 | 452 |
| **Test** $t2$ | 151 | 981 | 503 | 478 |
| **Total** | 1,683 | 10,480 | 5,355 | 5,125 |

Table 1: Overview ChartCheck statistics. Whilst claims are unique in each set, images may overlap except for $t2$ that contains images not present in any other split.

| Task | **R-**$\kappa$ |
|---|---|
| Chart filtering | 66.2 |
| Claim verification | 61.5 |
| Explanation verification | 51.4 |

Table 2: Randolph's Kappa (R-$\kappa$) as IAA scores.

caption. For example, workers identified and corrected partial explanations as well as errors in the explanation text. We created sets of seven tasks out of which two where gold-labelled, and included automated checks similar to Step 3. We calculated inter-annotator agreement (IAA) scores using Randolph's kappa for claim filtering, claim verification, and explanation verification tasks (see Table 2). Agreement over $61.0$ indicates substantial agreement while $51.4$ signifies moderate agreement (Landis and Koch, 1977).

### 3.3 Post-processing and validation (Step 5)

Finally, a group of annotators consisting of postgraduate computer science students evaluated all $1.9k$ testset claims and explanations from testsets $t1$ and $t2$ (see Table 1). We instructed them to evaluate and (if necessary) correct $(i)$ writing errors (e.g. typos, grammatical errors), $(ii)$ claim labels, $(iii)$ mismatch between claims and explanations. They also $(iv)$ adjusted "simple" claims, which require only one reasoning type for verification, to make them more complex.

## 4 Analysis of ChartCheck

### 4.1 Dataset Statistics

An overview of the ChartCheck dataset is given in Table 1. On average, the dataset contains six (three supporting and three refuting) claims per chart image. ChartCheck includes two testsets: while $t1$ is a testset has a similar distribution as the training set, for $t2$ we only included charts which are not part of the training set to evaluate models' performance on charts not seen during training.
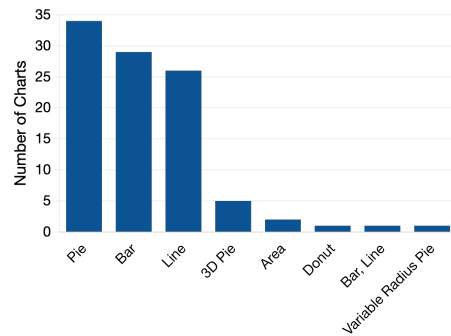


Figure 4: Chart types in annotated subset.

### 4.2 Chart Attributes

Following previous fact-checking datasets, e.g. Chen et al. (2020), we manually annotated a subset of 100 charts (i.e. 50 per testset) with the following attributes: chart type; complexity of charts; presence of datapoint labels, legends, axes titles, and main title; axes scales; background grid characteristics; number of charts per image; font size of title, labels, and other text occurring in the chart image; visual characteristics (e.g. colors and positions of legends). Fig. 4 shows the chart types found in the annotated subset. Around half of the charts ($47.5\%$) include labels their data points and $38.4\%$ have a legend describing the data categories. Most of these charts use colors ($32.3\%$) to group categories and only a small portion ($5.1\%$) uses different patterns. Charts that have an orientation (e.g. bar charts and line charts) are more often horizontal ($63.2\%$) than vertical ($36.8\%$). $7\%$ of images have multiple charts depicted in one image.

### 4.3 Reasoning with Charts

We labelled 230 dataset samples from both testsets with chart reasoning types based on a taxonomy of reasoning types humans use while interacting with chart data (Amar et al., 2005). Nine different types are present in the labelled subset: find extremum, comparison, world knowledge & commonsense (KCS), compute derived value, retrieve value, chart features, multichart reasoning, determine range, and multiaxis. Fig. 5 shows the reasoning types per testset. Around 100 out of all labelled claims require two to four reasoning types for correct classification of claims.

### 4.4 Model-generated Explanations

We further evaluated explanations by GPT4V (OpenAI, 2023), which is the best-performing model on claim classification, in detail. We asked postgradu-
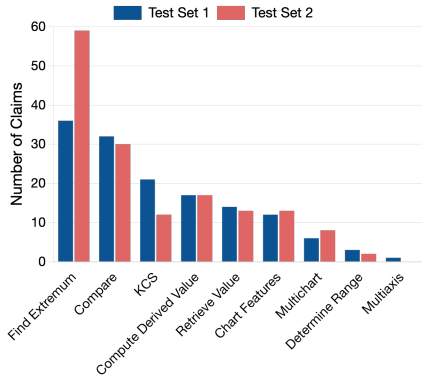
Figure 5: Reasoning types in annotated subset.

ate students to annotate 100 explanations based on the following five dimensions: $(i)$ factuality, $(ii)$ coverage, $(iii)$ readability, $(iv)$ coherence, and $(v)$ redundancy.

## 5 Experiments and Results

We define the explainable fact-verification task against chart evidence as follows. Each dataset entry $(s_i, v_i, c_i, y_i, e_i)$ comprises a natural language statement $s_i$, a chart image $v_i$, a caption $c_i$ accompanying the chart image, a verdict label $y_i \in \{0, 1\}$, and an explanation $e_i$ in natural language. The statement $s_i$ is a claim that is either supported $(y_i = 1)$ or refuted $(y_i = 0)$ by the chart (Fig. 1).

### 5.1 Baselines

We evaluate several baselines on ChartCheck (Fig. 6), which can be grouped into two categories: chart-to-table (CTT) and vision-language models (VLMs). The CTT baselines extract charts' underlying table data before performing classification and explanation generation. The VLMs are evaluated with Chain-of-Thought (CoT) prompting (Wei et al., 2022) in few- and zero-shot settings.

**Experimental setup.** Our baselines can be further grouped into single- and multi-task baselines. In the single-task setting, we first predict the claim verdict $\hat{y}_i$ and in a subsequent step generate the explanation using the claim, evidence, and either the predicted label $\hat{y}_i$ or gold label $y_i$ as input. In the multi-task setting, we train/instruct models to jointly *"classify and explain"*, resulting in output $(\hat{y}_i, \hat{e}_i)$. Inspired by the success of CoT on complex reasoning tasks, we also include the ChartCheck explanations in the input prompts as reasoning chains. We evaluate models in three different training settings: finetuning, few-, and zero-shot. The finetuned models are trained on the ChartCheck train-

ing set. For few-shot training, we provide three randomly selected training samples as input to the model.[8]

We evaluate claim classification using accuracy and macro $F_1$. We further evaluate explanations using BLEU (Papineni et al., 2002), ROUGE-L (Lin, 2004), BERTScore (Zhang et al., 2020), METEOR (Banerjee and Lavie, 2005), and BLEURT (Sellam et al., 2020).

(1) **Chart-to-table.** The chart-to-table architecture solves the task in a three-step approach. Inspired by Liu et al. (2022a), we decompose the chart-based fact-checking task into $(i)$ chart-to-table conversion, $(ii)$ table-based classification, and $(iii)$ explanation generation. First, for all charts we extract their underlying table data $T$ using the model DePlot (Liu et al., 2022a). We evaluate four models on fact verification with $(c_i, v_i, T_i)$ as input: DeBERTa (He et al., 2021), TAPAS (Herzig et al., 2020), FlanT5 (Shen et al., 2023), and GPT3.5.[9]

(2) **Vision-language models.** For our vision-language baseline, we use *MatCha* (Liu et al., 2022b) and GPT4-Vision (GPT4V) (OpenAI, 2023). MatCha is a finetuned version of the Pix2Struct model (Lee et al., 2022) for chart reasoning. It achieves SOTA performance on multiple chart datasets, including ChartQA (Masry et al., 2022b), PlotQA v2 (Methani et al., 2020), and the Pew charts subset of Chart-to-text (Kantharaj et al., 2022). We render a language prompt directly on top of the chart image. We evaluate three MatCha models which are pretrained on different chartQA datasets.

### 5.2 Results and Discussion

Table 3 and 4 provide an overview of baselines' on claim classification and explanation generation.

**What baseline performs best on ChartCheck? Claim Classification.** *DePlot-DeBERTa-class* achieves 75.0 accuracy on $t1$ and 72.5 on $t2$, out-performing other models. The baseline translates chart images to tables, concatenates the table text into a text sequence, and uses the claim, table, and the chart caption as input to a DeBERTa large model fine-tuned on ChartCheck. While translating chart to table data is beneficial on our task, not all claims are verifiable based on the extracted tables, for example, if the claim refers to colors visible in the image. Moreover, CTT models also propagate

---

[8] We provide the instruction text in Appendix Section D
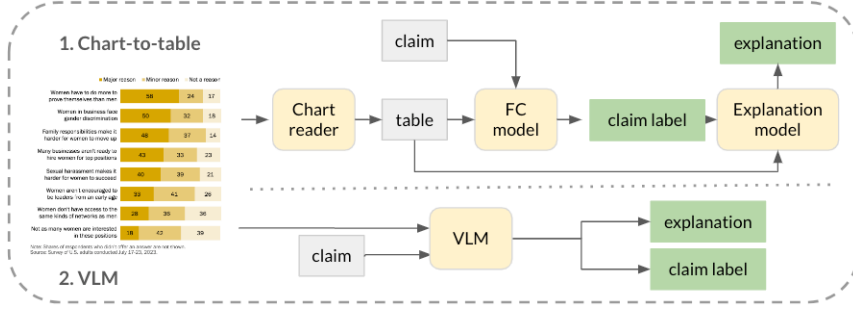[9] For further details we refer to Appendix D.

13925

Figure 6: Architecture of chart-to-table and vision-language model (VLM) baselines.

| Model | Task | Setting | Test1 Acc | Test1 $F_1$ | Test2 Acc | Test2 $F_1$ | Avg Test Acc |
|---|---|---|---|---|---|---|---|
| MatCha-chartqa-zero | C | zero | 29.1 | 46.4 | 29.3 | 48.3 | 29.2 |
| MatCha-plotqa1-zero | C | zero | 51.3 | 50.7 | 47.4 | 48.3 | 49.4 |
| MatCha-plotqa2-zero | C | zero | 48.2 | 32.9 | 48.7 | 32.9 | 48.4 |
| MatCha-finetune-class | C | fine-tuning | 64.0 | 63.7 | 61.6 | 60.9 | 62.8 |
| MatCha-finetune-multi | M | fine-tuning | 59.4 | 59.1 | 61.1 | 60.9 | 60.2 |
| MatCha-finetune-50/50-multi | M | fine-tuning | 61.9 | 61.2 | 63.0 | 62.1 | 62.4 |
| MatCha-chartqa-finetune-50/50-multi | M | fine-tuning | 59.1 | 58.8 | 61.4 | 60.7 | 60.2 |
| MatCha-plotqa1-finetune-50/50-multi | M | fine-tuning | 60.6 | 60.5 | 61.6 | 61.4 | 61.1 |
| MatCha-plotqa2-finetune-50/50-multi | M | fine-tuning | 62.8 | 62.8 | 61.4 | 61.4 | 62.1 |
| DePlot-DeBERTa-class | C | fine-tuning | **75.0** | **75.0** | **72.5** | **72.5** | **73.8** |
| DePlot-TAPAS-class | C | fine-tuning | 61.6 | 61.0 | 60.4 | 59.7 | 60.8 |
| DePlot-FlanT5-finetune-class | C | fine-tuning | 66.3 | 66.2 | 66.2 | 66.1 | 66.3 |
| DePlot-FlanT5-zero-multi | M | zero | 66.2 | 64.6 | 64.0 | 62.1 | 65.1 |
| DePlot-FlanT5-few-multi | M | few | 64.6 | 63.6 | 65.2 | 64.5 | 64.9 |
| DePlot-FlanT5-finetune-multi | M | fine-tuning | 65.7 | 65.7 | 65.9 | 65.8 | 65.8 |
| GPT3.5 | M | zero | 59.0 | 56.2 | 51.0 | 48.7 | 55.0 |
| GPT3.5 | M | few | 56.6 | 52.8 | 61.0 | 59.5 | 58.8 |
| GPT4V | M | zero | 73.8 | 73.5 | 72.0 | 71.3 | 72.9 |
| GPT4V | M | few | 74.0 | 73.6 | 70.6 | 70.2 | 72.3 |
| Human baseline | C | - | | | | | 95.7 |

Table 3: Classification results for models trained on claim classification (C) and in a multi-task setting with explanation generation (M).

errors to subsequent stages resulting in factually-incorrect explanations if the predicted claim label was wrong. Among the VLMs, the best-performing model is GPT4V that achieves 73.8. However, both models lag far behind human performance of 95.7.
**Explanation Generation.** On explanation generation, no baseline outperforms across all metrics. Overall, MatCha models underperform on explanation generation compared to GPT-models and CTT baselines. *DePlot-FlanT5-finetune-multi* outperforms other models on two out of five metrics. This model uses DePlot chart-to-table translation and finetunes FlanT5-base in a multitask setting on both claim classification and explanation generation. It performs best for testset $t1$ based on BERTScore (91.5) and ROUGE (46.3). For all other metrics, it scores among the top-3 baselines. We discuss the explanations' quality in more detail in Sec. 5.3.

**What model insights do we gain through different training settings? Pretraining.** Pretraining MatCha on further chart datasets (e.g. chartQA),

has no impact on the model's performance on ChartCheck. The redundancy of pretraining can be related to the datasets available for pretraining. Most of the questions in ChartQA/PlotQA only require retrieving a single value for the answer. ChartCheck claims' are complex and involve filtering, comparing, and commonsense knowledge among other reasoning types. Moreover, ChartCheck charts are more diverse (e.g. 3D pie charts and area charts) compared to other datasets.

**Multitask Learning.** Training models jointly on classification and explanation generation shows no positive effect for claim classification. While best-performing CTT and MatCha models are trained on classification-only, the multitask setting is more beneficial for explanation generation.

**Zero- and Few-shot Evaluation.** Comparing *DePlot-FlanT5-few-multi* to its zero-shot and fine-tuned equivalents, we find contrasting implications for claim classification and explanation generation. For claim classification, the accuracy for testset

| Model | Task | Training setting | BLEU | Rouge-L | MET | BERTScore | BLEURT |
|---|---|---|---|---|---|---|---|
| MatCha-finetune-multi | M | fine-tuning | 17.1 | 37.3 | 38.3 | 67.8 | 0.48 |
| MatCha-finetune-50/50-multi | M | fine-tuning | 17.3 | 37.1 | 37.9 | 67.3 | 0.45 |
| DePlot-FlanT5-zero-multi | M | zero | 0 | 5.04 | 1.3 | 83.2 | 0.48 |
| DePlot-FlanT5-few-multi | M | few | 5.2 | 28.5 | 21.1 | 90.6 | 0.56 |
| DePlot-FlanT5-finetune-multi | M | fine-tuning | 17.3 | **46.3** | 36.0 | **91.5** | 0.50 |
| DePlot-FlanT5-explanation-gold | E | fine-tuning | **30.7** | 38.7 | 33.2 | 89.6 | 0.57 |
| DePlot-FlanT5-explanation-pre-labels | E | fine-tuning | 12.7 | 37.9 | 32.4 | 89.5 | 0.60 |
| GPT3.5 | M | zero | 10.0 | 32.3 | **42.9** | 89.1 | 0.43 |
| GPT4V | M | few | 10.3 | 33.6 | 41.7 | 89.4 | 0.39 |

Table 4: Explanation generation results for testset $t1$ for model trained on explanation-only (E) and in a multi-task setting (M) with veracity classification.

| Chart Type | #Claims | DePlot | MatCha |
|---|---|---|---|
| Pie | 146 | 76.7 | 66.4 |
| Line | 122 | 62.2 | 59.0 |
| Bar | 104 | 81.7 | 64.4 |
| 3D Pie | 39 | 66.7 | 39.0 |
| Area | 5 | 60.0 | 40.0 |
| Donut | 2 | 50.0 | 50.0 |

Table 5: Performance of **DePlot**-DeBERTa-class and **MatCha**-finetune-class across different chart types.

$t1$ and $t2$ either does not change or decreases (e.g. from 65.9 to 64.0). However, prompting with few examples improves the model's explanation generation performance across all reported scores compared to its zero-shot counterpart.

## 5.3 Analysis with Performance Breakdown

**Which charts pose a challenge for the models?**
Comparing model performance across **chart types** (Table 5), we observe that the CTT baseline performs best on bar charts while MatCha excels at pie charts. While the performance on line and donut charts is comparable, the performance difference increases for pie, bar, and area charts. This can be related to certain chart types easier translatable to table data. CTT outperforms *MatCha-finetune-class* on single-chart images reaching 73.0 accuracy, whereas on **multi-chart images** the CTT baseline achieves only 50.0 accuracy, less than MatCha (58.3). Interestingly, both baselines struggle with charts with **legends**. The CTT baseline achieves 75.2 on charts without legends and 65.3 if a legend is given. Similarly, performance drops for the MatCha baseline from 65.2 (no legends) to 56.8 (legends).

**Do models struggle with specific types of chart reasoning?** While both baselines variants struggle with certain reasoning types, in general the best-performing chart-to-table (CTT) baseline outperforms the best MatCha model across all reasoning types (see Fig. 7). The difference is most significant for the reasoning types *determine range*,
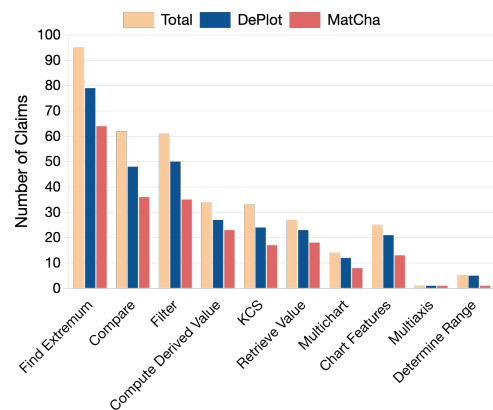


Figure 7: Number of correct predictions by reasoning types in labelled subset with **DePlot**-DeBERTa-class and **MatCha**-finetune-multi.
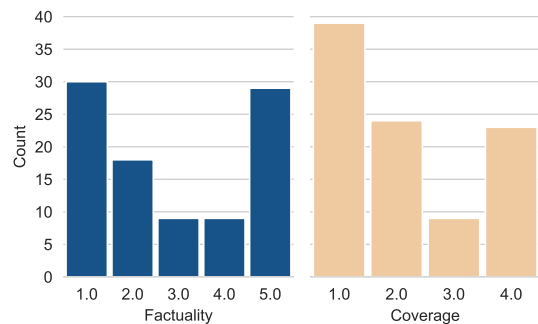


Figure 8: Manual rating of model-generated explanations for factuality (5 implies totally correct) and coverage (4 implies reference explanation fully covered).

*compare*, *chart features*, and *multichart*. CTT correctly predicts all claims requiring to determine a range of numbers correctly compared to the VLM reaching an accuracy of around 20.0. For the other three reasoning types, CTT correctly predicts most claims (approx. 80.0) while MatCha only achieves 60.0 or less correctly. Both models struggle with claims requiring commonsense reasoning (*KCS*).

**What are common failure cases in chart-to-table translation?** The charts in ChartCheck are extracted from the web and lack gold-tables required
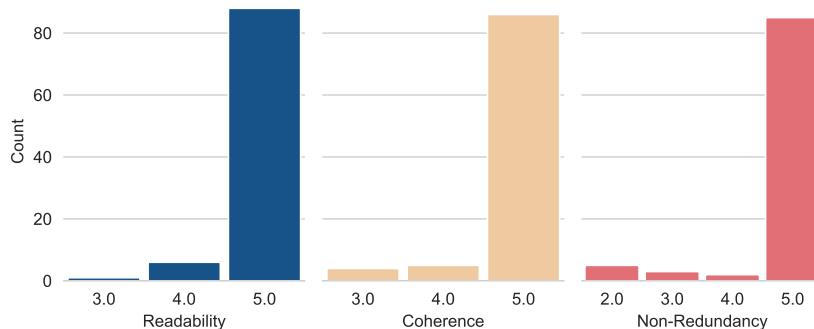
Figure 9: Results of manual explanation evaluation for the categories *Readability*, *Coherence*, *Non-Redundancy*. Expert annotators rated model-generated explanations on a scale from $1 - 5$ for these categories. Score 5 depicts *"very well readable"*, *"very coherent"*, and *"no redundant text"* for the given categories.

for automated evaluation. We manually evaluated a subset of DePlot-generated tables and found the following error cases:[10] (1) hallucinated labels for axes ticks and data points if they are not present or difficult to read; (2) errors in table translation of pie charts and grouped bar charts; (3) similar colors are confused if plotted against a non-white background; (4) missing or wrong titles; (5) incorrect text is extracted if background not white; (6) generated tables sometimes repeat one row multiple times or miss some data points.

**What are common failure cases in explanation generation?**  Rating model-generated explanations with the help of human annotators on factual correctness and coverage of the reference explanation, we find that approximately half of the explanations are either factually incorrect (i.e. equivalent to a score equal 1) or have major errors (score 2; see Fig. 8). A prevalent issue detected during assessment are explanations for refuted claims: models hallucinate explanations that contradict the content of the given chart but support the wrong claim. We also find explanations correctly mentioning the right claim label but giving a wrong explanation, ignoring the chart content. For example, a common issue is that models fail in reading numerical values from charts or in assigning them to the correct categories represented by a bar or a color in the image. Moreover, explanations include reasoning errors such as stating that $45\%$ is *"more than half"*. Expert annotators further rated explanations on a scale from $1 - 5$ for the categories *Readability*, *Coherence*, and *Non-Redundancy* (see Figures 9). More than $90\%$ of all explanations were rated *"very well readable"*, *"very coherent"*, and *"no redundant*

*text"*. Overall, these results imply that the model, while being able to generate very fluent, readable and coherent text, shows major limitations w.r.t. understanding and reasoning over charts to generate factually correct explanations.

## 6   Related Works

**Datasets for evidence-based fact checking**  The goal of evidence-based fact checking is to verify a claim against evidence data. Typically, evidence is conveyed in different forms, e.g. text (Thorne et al., 2018; Augenstein et al., 2019; Kotonya and Toni, 2020), tables (Aly et al., 2021; Akhtar et al., 2022; Chen et al., 2020), and images (Akhtar et al., 2023b). While recent work introduced a first chart fact-checking dataset, *ChartFC* (Akhtar et al., 2023a), the dataset has certain limitations. First, it is limited to bar charts, ignoring the diversity of data visualization found in real-world sources. Second, all ChartFC charts have been synthetically created using *Python* visualization libraries. Moreover, ChartFC claims are extracted from the Tab-Fact (Chen et al., 2020) dataset and therefore not written specifically for verification against charts.

**Chart datasets for other tasks**  Chart question answering (ChartQA) is a closely related task to ours. While existing ChartQA datasets have contributed to advancements in reasoning over data visualizations, most datasets have limitations in terms of size (Kim et al., 2021), template-based questions (Kafle et al., 2018), synthetically generated charts (Singh and Shekhar, 2020) or missing explanations. While ChartQA (Masry et al., 2022a) addresses many of these limitations, some challenges remain. First, ChartQA images are collected from four websites and while they repre-

---

[10]Examples are given in the appendix.

sent real-world scenarios, the charts are created by professionals, thus omitting the noise, design best-practise violations, and other issues found in charts by non-professional practitioners. Our dataset includes annotations for reasoning types and detailed visual attributes (e.g. presence of datapoint labels, legends, axis titles or background grid characteristics), which are absent in comparable datasets. Finally, ChartCheck is the first chart dataset with explanations, which has several benefits, such as better model training (i.e. allowing models to understand the reasoning behind a classification label), robustness, and transparency. This can further aid in building explainable models.

**Modeling approaches for charts**  Various approaches have been developed for chart modeling. One direction encodes both chart image and text, and applies a fusion method (e.g. attention) to combine the features into a joint representation (Masry et al., 2022b; Kafle et al., 2020). Another popular approach is to first extract the chart's table data, before applying table models (Kim et al., 2020; Liu et al., 2022a; Methani et al., 2020). Finally, more recent models use visual encoders to encode charts and their context and language decoders for text generation and solving NLP tasks, e.g. question answering (Masry et al., 2023; Ye et al., 2023; Liu et al., 2023; Xia et al., 2023). We evaluate ChartCheck on *DePlot* (Liu et al., 2022a), *MatCha* (Liu et al., 2022b), and *GPT4-Vision* (OpenAI, 2023) as SOTA models of each category.

## 7 Conclusion

We introduce ChartCheck, the first dataset for explainable fact-checking with $1.7k$ real-world chart and $10.5k$ human written claims and explanations. We evaluate SOTA models, including chart-to-table baselines and VLMs, in a few-/zero-shot and fine-tuned setting. Our best baseline yields 73.8 accuracy, lagging far behind human performance. Moreover, we study reasoning types and visual attributes that pose a challenge to SOTA models, as well as failures related to model-generated explanations.

## Limitations

The presented work exhibits the following limitations. First, the dataset is limited to English-only, whereas fact-checking misinformation is a global need that requires solutions in multiple languages. Further research is needed to extend chart-checking beyond English. Second, although ChartCheck includes a wide range of charts, certain chart types such as heat maps, scatter plots, and histograms are missing. Addressing these chart types should be a focus of future research. Finally, in our attempt to include diverse charts from real-world sources, we were unable to obtain gold table data for the charts. In the future, it would be beneficial to construct a dataset that includes diverse charts along with their underlying table data. This would enable the study of a broader range of models and their performance on such data. It would be also interesting to make the data three way labeled including "uncertain" label in addition to "true" and "false" claim.

## Ethics Statement

This paper introduces a dataset for automated fact-checking against chart images. The dataset is intended for research purposes, not as the sole means of evaluating real-world applications. The labels (i.e. supports and refutes) describe a claim's veracity given the evidence table. We do not make any statement on ChartCheck claims' truthfulness in a real-world context.

Prior to the crowdsourcing, we obtained ethical clearance from the relevant authority in our institution. We informed the annotators about all data being collected and its purpose. Participants had the opportunity to withdraw at any time and to provide feedback at the end of each task. All annotators were from English speaking countries. The payment was above the minimum wage and decided based on the time workers spent on the pilot tasks. In order to support future research, we plan to release all the scripts and resources used for dataset creation and model evaluation. This will facilitate and encourage further research in this field.

## Acknowledgements

## References

Mubashara Akhtar, Oana Cocarascu, and Elena Simperl. 2022. PubHealthTab: A public health table-based dataset for evidence-based fact checking. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1–16, Seattle, United States. Association for Computational Linguistics.

Mubashara Akhtar, Oana Cocarascu, and Elena Simperl. 2023a. Reading and reasoning over chart images for evidence-based automated fact-checking. In *Findings of the Association for Computational Lin-*

*guistics: EACL 2023*, pages 399–414, Dubrovnik, Croatia. Association for Computational Linguistics.

Mubashara Akhtar, Michael Schlichtkrull, Zhijiang Guo, Oana Cocarascu, Elena Simperl, and Andreas Vlachos. 2023b. Multimodal automated fact-checking: A survey. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5430–5448, Singapore. Association for Computational Linguistics.

Rami Aly, Zhijiang Guo, Michael Sejr Schlichtkrull, James Thorne, Andreas Vlachos, Christos Christodoulopoulos, Oana Cocarascu, and Arpit Mittal. 2021. FEVEROUS: fact extraction and verification over unstructured and structured information. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*.

Robert A. Amar, James Eagan, and John T. Stasko. 2005. Low-level components of analytic activity in information visualization. In *IEEE Symposium on Information Visualization (InfoVis 2005), 23-25 October 2005, Minneapolis, MN, USA*, pages 111–117. IEEE Computer Society.

Isabelle Augenstein, Christina Lioma, Dongsheng Wang, Lucas Chaves Lima, Casper Hansen, Christian Hansen, and Jakob Grue Simonsen. 2019. MultiFC: A real-world multi-domain dataset for evidence-based fact checking of claims. In *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 4685–4697, Hong Kong, China. Association for Computational Linguistics.

Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Sandra Carberry, Stephanie Elzer, and Seniz Demir. 2006. Information graphics: An untapped resource for digital libraries. In *Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '06, page 581–588, New York, NY, USA. Association for Computing Machinery.

Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyou Zhou, and William Yang Wang. 2020. Tabfact: A large-scale dataset for table-based fact verification. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking.

*Transactions of the Association for Computational Linguistics*, 10:178–206.

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: decoding-enhanced bert with disentangled attention. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.

Jeff Hemsley and Jaime Snyder. 2018. *FIVE Dimensions of Visual Misinformation in the Emerging Media Landscape*, pages 91–106. University of Texas Press, New York, USA.

Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Martin Eisenschlos. 2020. Tapas: Weakly supervised table parsing via pre-training. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 4320–4333. Association for Computational Linguistics.

Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrial-strength Natural Language Processing in Python.

Kushal Kafle, Brian L. Price, Scott Cohen, and Christopher Kanan. 2018. DVQA: understanding data visualizations via question answering. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 5648–5656. Computer Vision Foundation / IEEE Computer Society.

Kushal Kafle, Robik Shrestha, Brian L. Price, Scott Cohen, and Christopher Kanan. 2020. Answering questions about data visualizations using efficient bimodal fusion. In *IEEE Winter Conference on Applications of Computer Vision, WACV 2020, Snowmass Village, CO, USA, March 1-5, 2020*, pages 1487–1496. IEEE.

Shankar Kantharaj, Rixie Tiffany Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq Joty. 2022. Chart-to-text: A large-scale benchmark for chart summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4005–4023, Dublin, Ireland. Association for Computational Linguistics.

Dae Hyun Kim, Enamul Hoque, and Maneesh Agrawala. 2020. Answering questions about charts and generating visual explanations. In *CHI '20: CHI Conference on Human Factors in Computing Systems, Honolulu, HI, USA, April 25-30, 2020*, pages 1–13. ACM.

Wonjae Kim, Bokyung Son, and Ildoo Kim. 2021. Vilt: Vision-and-language transformer without convolution or region supervision. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 5583–5594. PMLR.

Neema Kotonya and Francesca Toni. 2020. Explainable automated fact-checking for public health claims. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7740–7754, Online. Association for Computational Linguistics.

J Richard Landis and Gary G Koch. 1977. The measurement of observer agreement for categorical data. *Biometrics*, pages 159–174.

Crystal Lee, Tanya Yang, Gabrielle D. Inchoco, Graham M. Jones, and Arvind Satyanarayan. 2021. Viral visualizations: How coronavirus skeptics use orthodox data practices to promote unorthodox science online. In *CHI '21: CHI Conference on Human Factors in Computing Systems, Virtual Event / Yokohama, Japan, May 8-13, 2021*, pages 607:1–607:18. ACM.

Kenton Lee, Mandar Joshi, Iulia Turc, Hexiang Hu, Fangyu Liu, Julian Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. 2022. Pix2struct: Screenshot parsing as pretraining for visual language understanding. *CoRR*, abs/2210.03347.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Maxim Lisnic, Cole Polychronis, Alexander Lex, and Marina Kogan. 2023. Misleading beyond visual tricks: How people actually lie with charts. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems, CHI 2023, Hamburg, Germany, April 23-28, 2023*, pages 817:1–817:21. ACM.

Fangyu Liu, Julian Martin Eisenschlos, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Wenhu Chen, Nigel Collier, and Yasemin Altun. 2022a. Deplot: One-shot visual language reasoning by plot-to-table translation. *CoRR*, abs/2212.10505.

Fangyu Liu, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Yasemin Altun, Nigel H. Collier, and Julian Martin Eisenschlos. 2022b. Matcha: Enhancing visual language pretraining with math reasoning and chart derendering. *CoRR*, abs/2212.09662.

Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen, Kaiqiang Song, Sangwoo Cho, Yaser Yacoob, and Dong Yu. 2023. MMC: advancing multimodal chart understanding with large-scale instruction tuning. *CoRR*, abs/2311.10774.

Leo Yu-Ho Lo, Yifan Cao, Leni Yang, and Huamin Qu. 2024. Why change my design: Explaining poorly constructed visualization designs with explorable explanations. *IEEE Trans. Vis. Comput. Graph.*, 30(1):955–964.

Leo Yu-Ho Lo, Ayush Gupta, Kento Shigyo, Aoyu Wu, Enrico Bertini, and Huamin Qu. 2022. Misinformed by visualization: What do we learn from misinformative visualizations? *Comput. Graph. Forum*, 41(3):515–525.

Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022a. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263–2279, Dublin, Ireland. Association for Computational Linguistics.

Ahmed Masry, Parsa Kavehzadeh, Xuan Long Do, Enamul Hoque, and Shafiq Joty. 2023. UniChart: A universal vision-language pretrained model for chart comprehension and reasoning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14662–14684, Singapore. Association for Computational Linguistics.

Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq R. Joty, and Enamul Hoque. 2022b. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 2263–2279. Association for Computational Linguistics.

Nitesh Methani, Pritha Ganguly, Mitesh M. Khapra, and Pratyush Kumar. 2020. Plotqa: Reasoning over scientific plots. In *IEEE Winter Conference on Applications of Computer Vision, WACV 2020, Snowmass Village, CO, USA, March 1-5, 2020*, pages 1516–1525. IEEE.

OpenAI. 2023. GPT-4 technical report. *CoRR*, abs/2303.08774.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.

Sheng Shen, Le Hou, Yanqi Zhou, Nan Du, Shayne Longpre, Jason Wei, Hyung Won Chung, Barret Zoph, William Fedus, Xinyun Chen, Tu Vu, Yuexin Wu, Wuyang Chen, Albert Webson, Yunxuan Li, Vincent Zhao, Hongkun Yu, Kurt Keutzer, Trevor Darrell, and Denny Zhou. 2023. Flan-moe: Scaling instruction-finetuned language models with sparse mixture of experts. *CoRR*, abs/2305.14705.

Hrituraj Singh and Sumit Shekhar. 2020. STL-CQA: Structure-based transformers with localization and

encoding for chart question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3275–3284, Online. Association for Computational Linguistics.

James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and verification. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 809–819. Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.

Renqiu Xia, Bo Zhang, Haoyang Peng, Ning Liao, Peng Ye, Botian Shi, Junchi Yan, and Yu Qiao. 2023. Structchart: Perception, structuring, reasoning for visual chart understanding. *CoRR*, abs/2309.11268.

Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Guohai Xu, Chenliang Li, Junfeng Tian, Qi Qian, Ji Zhang, Qin Jin, Liang He, Xin Lin, and Fei Huang. 2023. UReader: Universal OCR-free visually-situated language understanding with multimodal large language model. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2841–2858, Singapore. Association for Computational Linguistics.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with BERT. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

## A Future Challenges

Based on our experiments and results, we observed the following challenges as promising directions for future research: First, there is a need for models that can both extract (visual) information from



Figure 10: Chart distributed during the Brexit campaign[11].



Figure 11: Chart distributed during the Brexit campaign[12].

charts and successfully solve chart-related tasks that involve various reasoning skills. Although our DePlot-based baselines performed well, important visual information is lost during image-to-table translation. Second, there is a need for future datasets that expand on this work across various dimensions. One direction to explore is the creation of multilingual chart datasets. Another direction can consider fact-checking against data visualization types that are not currently included in ChartCheck, such as maps and infographics. Studying charts within the broader context of documents is another interesting direction to pursue. In real-world settings, charts are often embedded within documents. Evaluating claims about charts may require combining information from various sources, such as text, tables, or other images found alongside the chart.

## B Chart Misinformation

Fig. 10 and 11 are examples of chart-backed arguments spread during the Brexit referendum.

## C Dataset Collection

### C.1 Wikimedia Commons Repository

We collected charts using Wikimedia Commons, a project of the *Wikimedia Foundation* (simi-
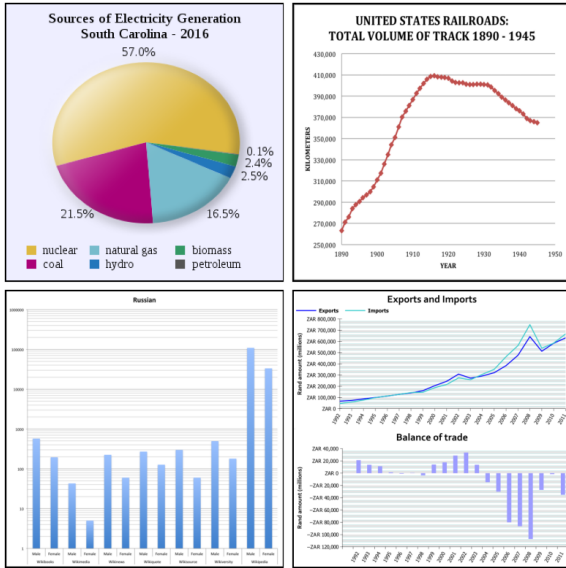
Figure 12: Examples of charts in the ChartCheck dataset.

lar to Wikipedia and WikiData).[13] Wikimedia is a media repository of images, videos, and audios which are used across all Wikimedia projects and not subject to copyright restrictions. As of May 2023, the repository contained over 93 million files. In Wikimedia Commons, images are labelled with categories which are part of a multi-hierarchy structure of categories. We used Wikimedia's API and the categories *"Horizontal_bar_charts"*, *"Vertical_bar_charts"*, *"Line_charts"*, *"Pie_charts_in_English"* (e.g. chart in Fig. 1), and *"Scatterplots"* to extract chart images from the repository. See Fig. 19 for a Wikimedia Commons entry. We only extracted charts, if their description was in English. We used the spaCy (Honnibal et al., 2020) language detection tool. Additionally, we extracted the file name, category, description, url, source (if available), and the Wikipedia pages the image is linked to.

## C.2 Crowdsourcing

**Mechanical Turk Details.** For the first and third crowdsourcing task, we created sets of seven annotation tasks each, out of which two were from a gold-labelled sample set annotated by the authors.

For the second annotation task (i.e. claim and explanation generation), we provided annotators a set of seven charts and captions per taskset. For each chart-caption pair, we asked them to write one claim which was supported by the chart and

Figure 13: Instructions we showed annotators at the start of the chart filtering task.

its caption and one which was refuted, as well as an explanation. The claim length was restricted to $[5; 30]$ tokens and the explanation text to $[10; 100]$ tokens. To improve the submission quality, we implemented automated checks on the claim and explanation text.

**Quality Assurance.** We included gold annotated samples in the first and final crowdsourcing task. We created a gold standard of $50$ charts for the chart filtering task and $40$ for the verification task. We used two gold samples per task. Workers who failed those two samples could not submit their work. For the second annotation task (i.e. claims and explanation generation), we manually evaluated one claim and explanation per submitted task set, before accepting the work. We banned workers with malicious behaviour, e.g. workers who submitted the same claims as supporting and refuting example. Automatic check evaluated claims and explanations for their token length, use of ambiguous wording (e.g. mostly, perhaps, and some), and negation words in case of refuting claims/explanations.

**Training and recruitment of workers.** We allocated the chart filtering and validation task to three crowdworkers each. Workers with a record of min. 1000 previously-approved tasks and an approval rate of min. $95\%$ were eligible. All workers first passed a chart literacy qualification test. We included examples of expert-labelled tasks in the instructions and rationales for the chosen labels.

Figure 13, 14, and 15 depict instructions presented to annotators at the beginning of the chart filtering, claim/explanation generation, and verification tasks, respectively.
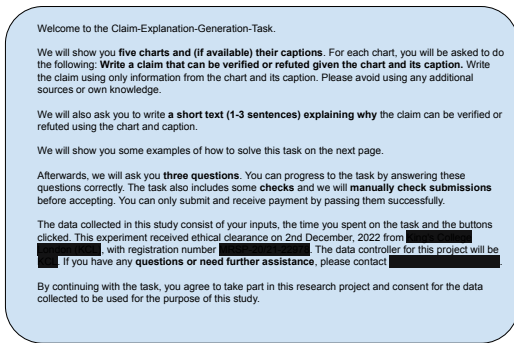
Figure 14: Instructions we showed annotators at the start of the claim/explanation generation task.
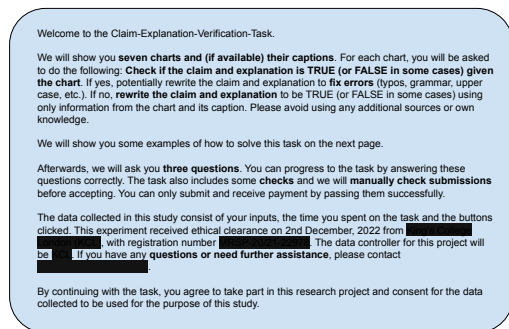


Figure 15: Instructions we showed annotators at the start of the claim/explanation verification task.

## C.3 Dataset Examples

Fig. 12 shows examples of charts in ChartCheck.

## D Modeling Details

**Technical details.** Before using images as input to models we converted a subset of images which were RGBA images to RGB images by adding a white chart background. We trained all MatCha models with AdaFactor optimizer, a learning rate of $1e-5$, a batch size of $4$ and a max patch size of 2048 for $15,000$ step. A checkpoint was saved every 300 step and the model with the best validation loss was selected. We used a max patch size of 2048. For DeBERTa-based classification, we evaluated and compared two models, the DeBERTa base model and a model pretrained on various NLI datasets (e.g. MNLI (Williams et al., 2018) and FEVER (Thorne et al., 2018)). We selected the NLI-pretrained DeBERTa model due to better performance and fine-tuned it with a learning rate of $1e-06$ and a batch size of $4$. We fine-tuned a TAPAS model which was previously trained on the table fact-checking task TabFact (Chen et al., 2020). The FlanT5 base model was fine-tuned on ChartCheck with a learning rate of $5e-5$ and batch

size of $8$.

**Instructions for Model Prompting.** For VLM experiments, our prompts consist of an instruction to the model with the claim (*"[claim]"*) and caption (*"[caption]"*) rendered on top of the input image. For claim classification tasks, we used the following prompt: *"The chart below has the following caption: [caption]. Given the caption and the chart below, is the following claim true: [claim]?"*. For claim classification with explanation generation task, we used the following prompt: *"The chart below has the following caption: [caption]. Given the caption and the chart below, is the following claim true: [claim]? Explain why?"*.

**Chart-to-table baseline details.** We finetune the first three models on the ChartCheck training data and evaluate GPT3.5 in a few- and zero-shot setting using CoT prompting. Finally, we finetune the FlanT5 (Shen et al., 2023) base model for explanation generation as described above. We also test FlanT5 in a multi-task setting (i.e. claim classification and explanation generation). We evaluate a finetuned, (*DePlot-FlanT5-finetune-multi*), few-shot, (*DePlot-FlanT5-few-multi*), and zero-shot (*DePlot-FlanT5-zero-multi*) version of the multi-task FlanT5 model.

**Vision-langauge model baseline details.** We evaluate three MatCha models in a zero shot setting on our testset: MatCha fine-tuned on ChartQA (*MatCha-chartqa-zero*), PlotQA v1 (*MatCha-plotqa1-zero*), and PlotQA v2 (*MatCha-plotqa2-zero*). In a fine-tuned setting, we have MatCha fine-tuning only on the classification task (*MatCha-finetune-classification*) and in a multi-task setting (*MatCha-finetune-multi*). Moreover, we experiment with a $50/50$ training data split during fine-tuning (*MatCha-finetune-50/50-multi*). For the first half of the dataset, we only fine-tune with classification prompts and for the second half we include explanations as well to improve overall classification performance.

## E Explanation Evaluation

Fig. 9 shows the results of manual explanation evaluation for the categories *Readability*, *Coherence*, *Non-Redundancy*. Expert annotators rated 100 model-generated explanations on a scale from $1-5$ for these categories. Score 5 depicts *"very well readable"*, *"very coherent"*, and *"no redundant text"* for the given categories.

**Chart:**



**Table:**
| TITLE | Frequency of mobile Wikipedia usage
Frequency of mobile Wikipedia usage | in | us
At least once a day | 29 | 13
At least once a month | 9 | 15
At least once a week | 36 | 45
Less than once a month | 6 | 10
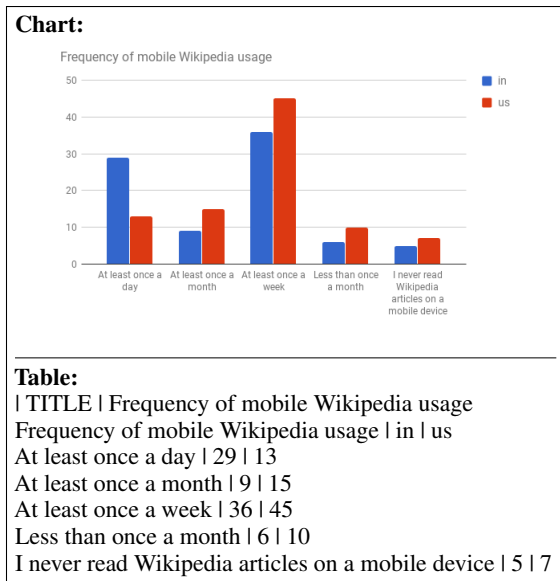I never read Wikipedia articles on a mobile device | 5 | 7

Figure 16: Example of DePlot Chart-to-Table conversion

## F  DePlot: Chart-to-table examples

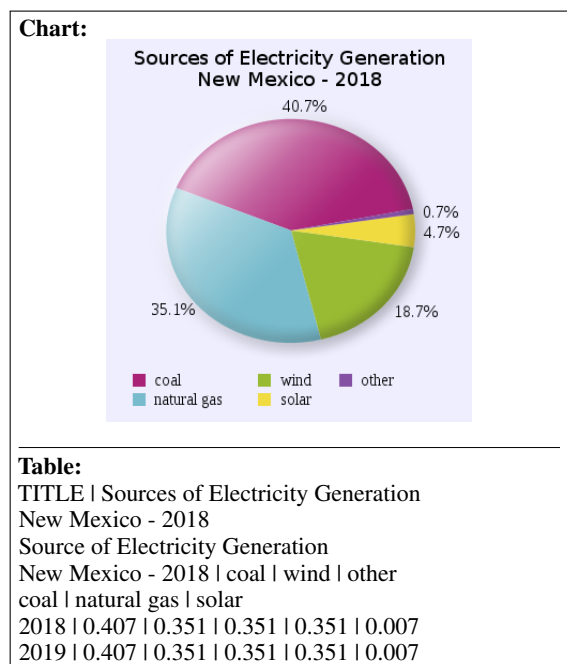Figures 16, 17, and 18 show ChartCheck images and extracted table data.

**Chart:**



**Table:**
TITLE | Sources of Electricity Generation
New Mexico - 2018
Source of Electricity Generation
New Mexico - 2018 | coal | wind | other
coal | natural gas | solar
2018 | 0.407 | 0.351 | 0.351 | 0.351 | 0.007
2019 | 0.407 | 0.351 | 0.351 | 0.351 | 0.007

Figure 17: Example of DePlot Chart-to-Table conversion

**Chart:**

United States Female Arson Arrests



**Table:**
TITLE | United States Female Arson Arrests
Age | Number arrested
9 | 0.36
10 | 0.82
13 | 2.74
15 | 3.11
16 | 1.83
17 | 2.01
18 | 1.27
19 | 1.44
20 | 1.31
21 | 1.52
22 | 1.62
23 | 1.80
24 | 1.67
25 | 1.69
30 | 1.39
35 | 1.35
40 | 1.33
45 | 0.99
50 | 0.72
55 | 0.52
60 | 0.23
65 | 0.11

Figure 18: Example of DePlot Chart-to-Table conversion

Figure 19: Chart example from the Wikimedia Commons page.