Scaling Text-Rich Image Understanding via Code-Guided Synthetic Multimodal Data Generation

Anonymous ACL submission

Abstract

Reasoning about images with rich text, such as charts and documents, is a critical application of vision-language models (VLMs). However, VLMs often struggle in these domains due to the scarcity of diverse text-rich visionlanguage data. To address this challenge, we present CoSyn, a framework that leverages the coding capabilities of text-only large language models (LLMs) to automatically create synthetic text-rich multimodal data. Given input text describing a target domain (e.g., "nutrition fact labels"), CoSyn prompts an LLM to generate code (Python, HTML, LaTeX, etc.) for rendering synthetic images. With the underlying code as textual representations of the 016 synthetic images, CoSyn can generate highquality instruction-tuning data, again relying on 017 a text-only LLM. Using CoSyn, we constructed a dataset comprising 400K images and 2.7M rows of vision-language instruction-tuning data. Comprehensive experiments on seven benchmarks demonstrate that models trained on our 022 synthetic data achieve state-of-the-art performance among competitive open-source models, including Llama 3.2, and surpass proprietary models such as GPT-4V and Gemini 1.5 Flash. Furthermore, CoSyn can produce synthetic pointing data, enabling VLMs to ground information within input images, showcasing its potential for developing multimodal agents capable of acting in real-world environments.

1 Introduction

033

037

041

Instruction-tuned vision-language models (VLMs) have shown strong performance across a range of multimodal tasks (Radford et al., 2021; OpenAI, 2023; Liu et al., 2023). However, these tasks typically focus on general image understanding over natural images rather than the specialized reasoning required for text-rich images such as charts, documents, diagrams, signs, labels, and screenshots. Understanding and reasoning over text-rich



Figure 1: Given a novel task (e.g., answering questions about nutrition facts), our code-guided generation system can produce targeted synthetic data to enhance the performance of VLMs on that specific task.

images is crucial for many applications, including analyzing scientific literature and figures (Asai et al., 2024), improving accessibility for users with visual impairments (Gurari et al., 2018), and enabling agentic workflows in real-world environments (Xie et al., 2024). Effectively interpreting these structured visual formats requires both textual comprehension and spatial reasoning, which current models struggle with due to the limited availability of high-quality, realistic, and diverse vision-language datasets (Methani et al., 2020).

To address these challenges and inspired by the fact that text-rich images are typically rendered from code, we develop **Code** Guided **Syn**thetic data generation system (**CoSyn**), a flexible framework for generating diverse synthetic text-rich multimodal data for vision-language instruction tuning. As illustrated in Figure 2, CoSyn can generate multimodal data for various target domains from a short natural language query, such as *book covers*. CoSyn leverages text-only LLMs, which excel at code generation, to produce both data and code that render diverse text-rich images using 11 supported rendering tools (e.g., Python, HTML, LaTeX). Grounded in the underlying code representation of the images, textual instructions are also 042

077

094

100

101

102

103

104

106

108

109

110

111

generated by the text-only LLM to create visionlanguage instruction-tuning datasets.

Using this framework, we construct the CoSyn-400K, as shown in Figure 3, a large-scale and diverse synthetic vision-language instruction-tuning dataset tailored for text-rich image understanding. We comprehensively evaluate the effectiveness of training on CoSyn-generated synthetic data across seven text-rich VQA benchmarks. Our model achieves state-of-the-art performance among competitive open-source models and surpasses proprietary models such as GPT-4V and Gemini 1.5. Notably, training on CoSyn synthetic data enables sample-efficient learning, achieving stronger performance with less data. In addition, CoSyn can synthesize chain-of-thought (CoT) reasoning data (Wei et al., 2022), improving performance on tasks requiring multi-hop reasoning. A fine-grained analysis of question types in ChartQA (Masry et al., 2022) reveals that training on CoSyn-400K results in stronger generalization to human-written questions. In contrast, models trained solely on existing academic datasets often overfit to biased training data, overperforming on templated or machinegenerated questions but struggling with more realistic, human-asked queries.

We then identify a key limitation of open-source VLMs that they struggle to generalize to out-ofdomain tasks they were not trained on. As shown in Figure 1, we introduce NutritionQA, a novel benchmark for understanding photos of nutrition labels, with practical applications like aiding users with visual impairments. Open-source VLMs perform poorly on this novel task, even after training on millions of images. However, by training on CoSyn-400K, our model adapts strongly to this novel domain in a zero-shot setting with significantly less training data. Remarkably, by generating just 7K in-domain synthetic nutrition label examples using CoSyn for fine-tuning, our model surpasses most open VLMs trained on millions of images. This highlights CoSyn's data efficiency and ability to help VLMs adapt to new domains through targeted synthetic data generation.

Finally, beyond the standard VQA task, we use 112 CoSyn to generate synthetic *pointing* training data, 113 which is particularly useful in agentic tasks. The 114 115 pointing data enables VLMs to retrieve coordinates for specific elements in a screenshot given a query 116 like "Point to the Checkout button" (Deitke et al., 117 2024). Our model trained on synthetic pointing 118 data achieves state-of-the-art performance on the 119

ScreenSpot click prediction benchmark (Baechler et al., 2024). Overall, our work demonstrates that synthetic data is a promising solution for advancing vision-language models in understanding text-rich images and unlocking their potential as multimodal digital assistants for real-world applications.

120

121

122

123

124

125

126

127

128

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

2 Related Work

Vision Language Models. Tsimpoukelli et al. (2021) first demonstrate that pre-trained, frozen language models can be extended to process visual inputs. Previous works fuse vision and language modalities using different strategies, such as cross-attention mechanisms (Alayrac et al., 2022) and Q-Former (Li et al., 2023). More recent architectures have converged on using MLP layers to project visual features into the language space (Liu et al., 2023). However, these architectures are often imbalanced, with the language backbone substantially larger than the visual encoder. As a result, without high-quality image-text data, models may overly rely on language priors, leading to hallucinations in their responses (Bai et al., 2024). Our work addresses this issue by generating highquality multimodal data for text-rich images.

Text-rich Images Understanding. Chart understanding and text-rich image understanding continue to challenge state-of-the-art models as naturally occurring vision-language data that can support training for understanding text-rich images is still scarce (Kahou et al., 2017; Kafle et al., 2018; Xu et al., 2023; Mukhopadhyay et al., 2024). In addition to charts and plots, a number of datasets address other kinds of text-rich images such as documents, infographics, diagrams and figures, and screenshots (Siegel et al., 2016; Mathew et al., 2021, 2022; Baechler et al., 2024; Roberts et al., 2024) have been made available. Many of these benchmarks are limited in size and scope, diversity of visualization types, and question types, making them suitable for evaluation but not for training data that could lead to generalized performance.

Synthetic Data for VLM. Generating synthetic images with annotations grounded in known source representations has been widely used in domains with limited vision-language data (Johnson-Roberson et al., 2017; Johnson et al., 2017; Cascante-Bonilla et al., 2022; Zhang et al., 2024). This approach has been applied to chart and plot VQA typically using a limited small set of chart types and by instantiating handcrafted question



Figure 2: The overview of our **Co**de Guided **Syn**thetic data generation system (**CoSyn**), which has 20 generation pipelines based on 11 render tools. Given a user query, e.g., "book cover," CoSyn selects the appropriate pipelines and starts with generating diverse topics conditioned on personas, then synthesizes detailed data for code generation. The code renders the image and is also fed as context for an LLM to construct instruction-tuning data.

templates (Kahou et al., 2017; Kafle et al., 2018; Methani et al., 2020; Singh and Shekhar, 2020). Following this, Li and Tajbakhsh (2023) and Carbune et al. (2024a) explore using text-only LLMs to generate annotations or Q&A pairs from table or text descriptions associated with charts to train VLMs. Other recent approaches, similar to the procedure in this work, explore generating data and code to render synthetic charts (Han et al., 2023; Shinoda et al., 2024; Xia et al., 2024) while using the data and code representation to generate annotations and Q&A. These works generate synthetic data that is still highly limited in terms of the diversity of topics, figure types, and rendering pipelines, which is important for generalizing to out-of-distribution tasks. In our work, we expand the scope of our generation beyond charts to encompass a wider range of diverse text-rich images.

170

171

172

173

174

175

177

178

179

182

183

184

186

187

190

192

193

194

196

198

201

204

3 Problem Formulation

Given a text query q about an image type, e.g., flow charts, our goal is to create a synthetic multimodal dataset $\mathcal{D}_q = \{(I,T)\}$, where I is the image, and T is the textual instruction-tuning data (e.g., question-answer pairs). \mathcal{D}_q is used to train a VLM to improve its ability to understand images related to q. The core idea of our approach is using code C as the intermediate representation to bridge the image and text. The overall generation process can be decomposed as follows:

$$P(I,T|q) = P_{LM}(C|q) \cdot P(I|C) \cdot P_{LM}(T|C)$$

where $P_{\text{LM}}(C|q)$ represents prompting a language model to generate code C, which is executed to render the image, P(I|C). $P_{\text{LM}}(T|C)$ uses code C (without the image) as context for an LLM to generate the textual instruction-tuning data.

4 CoSyn System

Figure 2 illustrates the workflow of our **Code**-Guided **Syn**thetic data generation system (**CoSyn**). The system takes a language input, such as "generate a dataset of book covers", and outputs a multimodal dataset. Based on the input query, CoSyn selects one of 20 generation pipelines built on 11 rendering tools. The process starts with topic generation, conditioned on a sampled persona that guides the style and content. Next, the system generates data content and converts it into code, which is then executed to render synthetic images. Finally, using the code as context, we prompt the LLM to generate corresponding textual instructions.

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

In the following, we provide detailed explanations of the rendering tools supported by CoSyn, the tailored generation pipelines based on these tools, our persona-driven approach to diversify content and styles, and the large-scale dataset of 400K synthetic images generated by CoSyn.

Rendering Tools. We integrate various rendering tools to generate diverse types of images, forming the foundation of CoSyn's ability for text-rich image generation. For example, Matplotlib, Plotly, and Vega-Lite are used to create different types of charts. LaTeX and HTML are used for documents and tables, while Mermaid and Graphviz generate diagrams. We utilize SVG and Asymptote to create vector graphics and math-related content. For specialized tasks, we rely on Lilypond to generate music sheets and RDKit for chemical structures. We implement customized functions for each tool to execute LLM-generated code and obtain corresponding rendered images. These tools collectively enable CoSyn to produce a wide range of highquality, text-rich synthetic images.

Pipelines. We design 20 pipelines based on 11



Figure 3: Our CoSyn-400K dataset consists of 9 categories of text-rich images with 2.7M instruction-tuning data. More qualitative examples, along with question-answer annotations, are available in Figure 12 -18 in Appendix C.

rendering tools.¹ Each pipeline follows the same procedure: (1) *Topic generation* to define the theme of this synthetic example, (2) *Data generation* to populate the detailed contents, (3) *Code generation* to create executable code that renders the image, and (4) *Instruction generation* uses the code as context to produce textual instruction-tuning data. Each stage is controlled by tailored prompts customized for the image category and rendering tool. Figure 8 shows an example of all prompt templates used in the HTML Document pipeline.

242

247

248

249

250

261

265

Use personas for diversity. LLMs often struggle to generate diverse synthetic data using sampling parameters alone (Yu et al., 2023), with biases leading to repetitive outputs across different runs. Recent work (Ge et al., 2024) shows that incorporating personas in prompts can improve diversity by enabling models to generate from varied perspectives. CoSyn adopts personas to enhance diversity during the Topic Generation stage. Each persona is a short sentence describing a personality or identity. For example, as shown in the middle of Figure 2, we sample a persona "a sci-fi novelist who likes alien worlds", which results in a topic of "a novel about Extraterrestrial Flora & Fauna" for generating the book cover image. We use the 200K personas released by Ge et al. (2024).

269Implementation details. CoSyn is built on the270DataDreamer library (Patel et al., 2024), which271supports robust multi-stage synthetic data genera-272tion pipelines that are easy to maintain, reproduce,273and extend. DataDreamer documents the prompts274and parameters used at each generation stage and

implements several efficient techniques, such as parallel generation and response caching, to optimize performance. For the data and code generation stages, we use Claude-3.5-Sonnet, which performs well in coding tasks (Anthropic, 2024b). For instruction-tuning data generation, we select GPT-4o-mini (OpenAI, 2023) for its cost efficiency. 275

276

277

278

279

280

281

284

285

287

288

290

291

292

293

294

296

297

298

299

300

301

302

303

304

305

307

CoSyn-400K. As shown in Figure 3, we use CoSyn to generate a large-scale synthetic dataset of 400K images across nine categories: charts, documents, math problems, tables, diagrams, vector graphics, music sheets, electrical circuits, and chemical structures. Since CoSyn is controlled via language inputs, it can easily generate diverse, fine-grained image types by varying the input queries. For instance, we use over 100 queries to generate document data covering *receipts, resumes, meal plans*, etc. Some queries used for CoSyn-400K are provided in Appendix A.3. This ensures that our dataset covers a broad range of domains. The following sections validate how our synthetic datasets enhance the ability of VLMs to understand text-rich images.

5 Experimental Setup

Our experiments aim to verify the value of our synthetic data in the supervised fine-tuning stage of training vision-language models. This section introduces the architecture of our model, training strategy, datasets we used, baselines for comparison, and other details on implementation.

Model Architecture. We follow the same image preprocessing and architecture as Molmo (Deitke et al., 2024), which uses the MLP layer to connect the vision encoder and a pretrained LLM. We

¹Some tools are used in multiple pipelines, e.g., HTML is used for generating documents, tables, and charts.

Model	ChartQA	DocVQA	InfoVQA	TableVQA	AI2D	TextVQA	ScreenQA	Average
GPT-4V	78.1	87.2	75.1	60.5	89.4	78.0	41.6	72.8
Gemini 1.5 Flash	85.4	89.9	75.3	72.6	91.7	78.7	40.1	76.2
Claude-3 Opus	80.8	89.3	55.6	70.0	88.1	67.5	39.8	70.2
PaliGemma-3B [†]	71.4	84.8	47.8	46.4	73.3	76.5	32.2	61.8
BLIP-3-4B [†]	60.0	61.4	31.5	24.3	74.2	71.0	26.2	49.8
Cambrian-7B [†]	73.3	77.8	41.6	40.6	73.0	71.7	44.4	64.2
LLaVA-1.5-7B [†] *	17.8	28.1	25.8	33.1	55.5	58.2	17.6	33.7
LLaVA-Next-8B [†]	69.5	78.2	43.8	43.9	71.6	65.3	34.2	58.1
LLaVA-OneVision-7B [†]	80.0	87.5	68.8	64.6	81.4	78.3	46.3	72.4
Pixtral-12B	81.8	90.7	50.8	67.0	79.0	75.7	39.4	69.2
Llama 3.2 11B	83.4	88.4	63.6	51.1	91.9	73.1	87.7	<u>77.0</u>
Ours (7B) [†]	86.3	90.0	70.5	65.8	91.9	82.0	80.1	80.9
Ours (7B-zero-shot) ^{†*}	80.8	82.9	59.8	64.9	83.9	72.7	78.1	74.7

Table 1: **Results on 7 text-rich benchmarks.** The result of the best-performing open-source model is **bold**, and the second-best is <u>underlined</u>. Models with \dagger stand for open data and code for multimodal training. Models with * are zero-shot models, which means the models are not trained on instances from any of the evaluation datasets.

choose OpenAI's CLIP (ViT-L/14 336px) (Radford et al., 2021) as the vision backbone and Mistral-7B (Jiang et al., 2023) as the language model.

310

311

312

313

314

315

316

317

319

322

327

328

329

334

335

339

340

Training Process. We adopt the same training strategy as Molmo (Deitke et al., 2024), which consists of two stages: (1) *Pre-training* on dense captions from PixMo-Cap and (2) *Supervised fine-tuning* on three categories of datasets below:

Evaluation Datasets. We evaluate our model on seven text-rich benchmarks, including ChartQA (Masry et al., 2022), DocVQA (Mathew et al., 2021), InfographicVQA (Mathew et al., 2022), TableVQA-Bench (Kim et al., 2024), AI2 Diagrams (Kembhavi et al., 2016), TextVQA (Singh et al., 2019), and ScreenQA (Baechler et al., 2024). We adopt their official metrics for calculating performance. In total, we have 138K training images from the evaluation datasets.²

- Auxiliary Datasets. We select additional academic datasets for fine-tuning: VQAv2 (Goyal et al., 2017), GQA (Hudson and Manning, 2019), OK-VQA (Marino et al., 2019), OCR-VQA (Mishra et al., 2019), A-OKVQA (Schwenk et al., 2022), ScienceQA (Lu et al., 2022), TabMWP (Lu et al., 2023), ST-VQA (Biten et al., 2019), TallyQA (Acharya et al., 2019), DVQA (Kafle et al., 2018), FigureQA (Kahou et al., 2017), and PlotQA (Methani et al., 2020). The auxiliary datasets contain around 1M training images.
- Synthetic Datasets. As introduced in Sec 4 and also shown in Figure 3, our synthetic datasets include 400K text-rich images from 9 categories.
 Our best-performing model uses all three cate-

gories of datasets above. We also trained a zeroshot model using only auxiliary and synthetic data without any examples from the evaluation datasets, which still exhibits competitive benchmark performance, as shown in the last row of Table 1. 341

342

343

344

345

347

348

349

351

352

353

354

355

356

357

358

360

361

362

363

364

365

366

367

369

370

371

372

373

Baselines. We compare our model against recent open-source VLMs with a similar scale (7B), including PaliGemma-3B (Beyer et al., 2024), BLIP-3-4B (Xue et al., 2024), Cambrian-7B (Tong et al., 2024), LLaVA-1.5-7B (Liu et al., 2023), LLaVA-Next-8B (Liu et al., 2024), LLaVA OneVision-7B (Li et al., 2024), Pixtral-12B (Agrawal et al., 2024). We also compare with proprietary models: GPT-4V (OpenAI, 2023), Gemini-1.5-Flash (Team, 2024), and Claude-3 Opus (Anthropic, 2024a).

Implementation Details. We train our model on TPU v3-128 with a batch size of 32. Our best-performing model is trained for 60K steps, taking about 30 hours. The checkpoints with the highest validation performance are retained for testing.

6 Results

This section covers (1) the competitive performance of the model trained on our synthetic data (Sec 6.1), (2) the comprehensive analyses to highlight the benefits of synthetic data (Sec 6.2), and (3) the effectiveness of synthetic pointing data in improving VLMs for web agent tasks (Sec 6.3).

6.1 Main Results

Table 1 compares our model's performance with both open and closed models across seven text-rich benchmarks. On average, our 7B model achieves the highest performance, surpassing the secondbest model (Llama 3.2 11B) by 3.9%. Notably, our

²TableVQA is an eval-only benchmark (no training split), and we do not use the training split from ScreenQA.



Figure 4: Ablation on training data selection. Aux, Syn, and Eval stand for auxiliary, synthetic, and evaluation datasets, respectively. We report the average score on eight benchmarks. The detailed performance breakdown on each benchmark is in Table 7.

model ranks first in four out of the seven datasets and second in the remaining three. More surprisingly, our zero-shot model (the last row in Table 1) outperforms most open and closed models without exposure to any training instances from the evaluation datasets. In contrast, these competing models often rely on benchmark training data and are thus not true zero-shot models. This result demonstrates that the capabilities learned from our synthetic data can transfer effectively to downstream tasks.

Analysis 6.2

In the following experiments, we quantify the contribution of synthetic data to the benchmark performance by ablating the combinations of fine-tuning datasets. Then, we demonstrate that our CoSyn system can efficiently assist VLMs in generalizing to novel tasks. Finally, we show that synthetic data can help mitigate the overfitting of biases.

Synthetic data boosts the performance. Table 4 presents an ablation study on the choices of supervised fine-tuning data. In the zero-shot settings, when the model is trained on auxiliary datasets (over 1M training images not directly from the evaluation tasks), it fails to generalize effectively to the evaluation tasks, with a substantial performance gap of 14.1% below GPT-4V. However, using only 400K synthetic samples achieves a perfor-400 mance comparable to GPT-4V. Our best zero-shot model surpasses GPT-4V when jointly training syn-402 thetic and auxiliary data. Under the supervised 403 settings, training with in-domain data alone yields 404 405 strong performance. However, adding 1M auxiliary samples provides a modest improvement of 1.4%, 406 while incorporating synthetic data results in a more 407 significant 3.6% boost. These results demonstrate 408 the effectiveness of synthetic data in enhancing 409



Figure 5: Zero shot performance on NutritionQA. The x-axis denotes the number of training examples used for the instruction-tuning stage. The models on the upper left side demonstrate better data efficiency.

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

VLMs' performance on text-rich tasks.

Zero-shot Generalization on a Novel Task. Vision-language models typically rely on indomain data to perform well on specific tasks. When encountering a novel task, such as answering questions about nutrition labels in Figure 1, models without seeing similar examples during training may struggle with this novel task. However, our CoSyn system enables controllable data generation. Given the task name as input, CoSyn can generate task-specific data to fine-tune the model.

To validate this, we annotated a small evaluation dataset called NutritionOA, which includes 100 examples of questions about photos of nutrition labels. Some questions require multi-hop reasoning, as Figure 10 illustrates. We evaluated GPT-4V and several open-source VLMs on this dataset and report the performance in Figure 5. The x-axis in Figure 5 represents the amount of data used during the instruction fine-tuning stage.

Despite being trained on millions of images, we observe that open-source VLMs are not dataefficient and perform poorly on this novel task compared to GPT-4V. Although many open-source VLMs claim to achieve GPT-4V-level performance, they fall short when tested on new tasks in the wild. Without synthetic data, our model (Eval + Aux) achieves results similar to those of open models. However, when trained on 400K synthetic samples, our model matches GPT-4V's performance.

More impressively, we used CoSyn to generate 7K synthetic nutrition label samples and fine-tuned the model using only this 7K data. The resulting model outperforms most open-source VLMs on the NutritionQA task. These results demonstrate that code-guided synthetic data is an effective and efficient method for adapting VLMs to new domains.

401

374



Figure 6: Ablation of using Chain-of-Thought reasoning. Short Answer represents prompting model to output the answer as short as possible. + CoT stands for providing Chain-of-Thought reasoning before giving the final answer. Results on all datasets are in Table 6.

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

Synthetic Data for Chain-of-Thought Reasoning. Existing text-rich datasets, such as ChartQA (Masry et al., 2022), are typically annotated with short answers. However, questions like "Compute the mean of the data in the plot" require step-bystep mathematical reasoning to arrive at the correct answer. Models trained only with short-answer supervision may fail to learn proper plot comprehension, but instead overfitting to annotation biases in these datasets. On the contrary, our CoSyn-400K includes explanation text alongside the short answer. Each instruction-tuning example consists of a (question, explanation, short answer) triplet, enabling models to learn chain-of-thought (CoT) reasoning. During fine-tuning, we design two prompt templates for our synthetic data:

CoT Prompt: **<Question>** Provide reasoning steps and then give the short answer. **<Explanation>** Answer: **<Answer>**

Short Answer Prompt: **<Question>** Answer with as few words as possible. **<Answer>**

Those prompts allow VLMs to switch between the two answering styles and perform CoT reasoning when necessary. Figure 6 shows that incorporating CoT reasoning improves performance on ChartQA, TableVQA, and NutritionQA, as these datasets contain examples requiring multi-hop reasoning. However, we observe that adding CoT reasoning reduces performance on DocVQA and InfoVQA. We find this decline is caused by answer biases in these benchmarks. Specifically, the ground-truth answers favor short responses, often penalizing more detailed and verbal responses. For instance, in DocVQA, the ground-truth for an example is "T-Th", whereas the model responds with "Tuesday to Thursday". Although the response is



Table 2: **Results on human and machine-generated questions of ChartQA.** The pie charts above display the percentage distribution of two question types in training and testing. Δ (\downarrow lower is better) denotes the performance gap between human and machine questions.

correct, the strict string-matching metric assigns it a zero score. This highlights key limitations of current multimodal benchmarks, including answering biases and rigid evaluation metrics that fail to capture the full extent of a model's capabilities. 480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

502

503

505

506

507

508

509

510

511

512

513

Synthetic Data for Mitigating Biases. Our previous experiments reveal answering biases in multimodal benchmarks, which VLMs trained solely on these datasets often inherit. To further validate this issue, we analyze ChartQA and observe a distribution shift in question types. As shown in the pie charts above Table 2, some ChartQA questions are human-annotated, while others are generated by the language model T5 (Raffel et al., 2020), which is heavily influenced by prompt phrasing and limited to a fixed set of question templates. During training, most questions (73.9%) in ChartQA are machine-generated, while the test set contains an even distribution of human-annotated and machinegenerated questions. Models trained exclusively on ChartQA tend to overfit to T5-generated questions. Table 2 illustrates this issue: PaliGemma (Beyer et al., 2024) and ChartPali (Carbune et al., 2024b) achieve high accuracy on machine-generated questions but experience a significant performance drop of over 30% on human-annotated questions.

Similarly, without synthetic data, our model shows a noticeable 21.8% gap between the two question types. However, incorporating synthetic data during training reduces this gap to 14.2%, improving the model's ability to answer human-asked questions. This suggests that synthetic data can mitigate overfitting on benchmarks and enhance VLMs' usability in real-world applications.



(a) Synthetic Pointing Data Generation

(b) Real Screenshot Pointing

Figure 7: **The overview of enabling VLMs to point through synthetic data.** (a) We synthesize pointing data by prompting an LLM to generate pointing questions and edit the code to draw the answer points explicitly. (b) We demonstrate that the VLM trained on synthetic pointing data can be generalized to real agentic tasks.

6.3 Synthetic Pointing Data

514

515 516

517

518

519

520

521

Pointing enables vision-language models to answer questions by providing specific points on images. This functionality allows models to ground their responses in visual content and interact with environments, which is crucial for developing digital agents. We find that we can synthesize pointing data using our code-guided generation system.

Method. Since we have access to the source code 522 for all generated images, we can prompt an LLM 523 to modify the code to draw points on the images 524 explicitly. As illustrated in Figure 7, we feed the 525 526 image's source code as context to the LLM, which generates a pointing question and edits the code to draw points with a predefined color. By extracting the pixel values of these points, we can obtain their exact (x, y) coordinates.³ We then use this point-530 ing data to train VLMs, enabling them to answer questions by providing point coordinates. In total, 532 we generate pointing data for 65K synthetic images. Figure 19 shows some qualitative examples from our synthetic pointing dataset.

Setup. We evaluate pointing ability on ScreenSpot 536 (Cheng et al., 2024), where the task requires models to provide the correct click location based on a 538 given instruction. ScreenSpot contains screenshots 539 from mobile phones, desktops, and web pages. To 540 assess the effectiveness of our synthetic pointing data, we compare it to the model trained on PixMo-542 point (Deitke et al., 2024), which consists of 155K 543 human-annotated images. Our best-performing 544 model uses both PixMo-point and synthetic pointing data. Additionally, we compare against existing methods like CogAgent (Hong et al., 2024), 547 SeeClick (Cheng et al., 2024), and UGround (Gou 548 et al., 2024), which is trained on 1.3M screenshots. 549

	Мо	bile	Des	ktop	W	eb	Avg
Model	Text	Icon	Text	Icon	Text	Icon	
GPT-40	20.2	24.9	21.1	23.6	12.2	7.8	18.3
CogAgent	67.0	24.0	74.2	20.0	70.4	28.6	47.4
SeeClick	78.0	52.0	72.2	30.0	55.7	32.5	53.4
UGround	82.8	<u>60.3</u>	82.5	<u>63.6</u>	<u>80.4</u>	70.4	73.3
Synthetic	90.8	53.3	78.4	58.6	80.0	47.1	68.0
Human	84.2	59.0	88.1	52.9	76.5	50.5	68.5
Combined	<u>89.0</u>	65.1	<u>87.6</u>	65.7	83.0	<u>58.7</u>	74.9

Table 3: **Click accuracy on ScreenSpot.** We report our models trained on different pointing data. Human stands for using the human-annotated data from PixMopoint (Deitke et al., 2024). Combined means combining human-annotated data with our synthetic pointing data.

550

551

552

553

554

555

556

557

558

559

560

561

562

564

565

566

567

568

570

571

572

Results. Table 3 compares the click accuracy of our models with previous methods. Using 65K synthetic pointing samples, our model achieves performance comparable to the one trained on 155K human-annotated samples. When combining synthetic and human data, our model achieves state-of-the-art performance on ScreenSpot, surpassing the recent UGround (Gou et al., 2024), which was trained on 1.3M screenshots. These results demonstrate that synthetic pointing data is a data-efficient approach for improving VLM performance on agentic tasks involving click prediction.

7 Conclusion

In this work, we introduced CoSyn, a framework for generating synthetic data that significantly enhances VLM performance on text-rich image understanding. Our comprehensive analysis highlights the advantages of synthetic data for domain generalization, data efficiency, and bias mitigation. Our work demonstrates that the coding capabilities of text-only LLMs can effectively assist multimodal learning and unleash the potential of visionlanguage models for real-world applications.

 $^{^{3}}$ The coordinates are normalized to (0, 100) to mitigate the influence of image resolution.

Limitation 573

The effectiveness of synthetic data depends heavily on the quality and diversity of the prompts and 575 rendering pipelines used for data generation. For highly specialized or underrepresented domains, generating sufficiently diverse data remains challenging and may require careful prompt engineering or additional customization of rendering tools. 580 Targeted synthetic data generation may be essential 581 for certain tasks to achieve adequate performance, and ensuring relevance and coverage still requires 583 domain-specific expertise. Synthetic data also may 584 not fully capture the complexity of real-world data 585 in some scenarios. Therefore, improving the diversity and realism of synthetic data to better support 587 models in highly variable or evolving domains is a 588 reasonable avenue for future research. Finally, our current synthetic data is limited to English and may require further extension for multilingual support. 591

Ethical Statement

To the best of our knowledge, this work presents 593 no significant ethical concerns. We note, however, 595 that the use of synthetic data can propagate biases present in the generation model used. Conversely, 596 synthetic data can also help mitigate biases and 597 expand coverage, as demonstrated in this work, by 598 greatly expanding the domains present in visionlanguage instruction-tuning training data to yield stronger generalized performance.

References

602

605

611

614

615

617

- Manoj Acharya, Kushal Kafle, and Christopher Kanan. 2019. TallyQA: Answering complex counting questions. In AAAI.
- Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Devendra Chaplot, Jessica Chudnovsky, Saurabh Garg, Theophile Gervet, Soham Ghosh, Amélie Héliou, Paul Jacob, et al. 2024. Pixtral 12b. arXiv preprint arXiv:2410.07073.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. In NeurIPS.
- Anthropic. 2024a. The claude 3 model family: Opus, 616 sonnet, haiku.
- Anthropic. 2024b. Introducing the next generation of 619 claude.

Akari Asai, Jacqueline He, Rulin Shao, Weijia Shi, Amanpreet Singh, Joseph Chee Chang, Kyle Lo, Luca Soldaini, Sergey Feldman, Mike D'arcy, et al. 2024. Openscholar: Synthesizing scientific literature with retrieval-augmented lms. arXiv preprint arXiv:2411.14199.

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

- Gilles Baechler, Srinivas Sunkara, Maria Wang, Fedir Zubach, Hassan Mansoor, Vincent Etter, Victor Cărbune, Jason Lin, Jindong Chen, and Abhanshu Sharma. 2024. Screenai: A vision-language model for ui and infographics understanding. Preprint, arXiv:2402.04615.
- Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou. 2024. Hallucination of multimodal large language models: A survey. arXiv preprint arXiv:2404.18930.
- Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, Thomas Unterthiner, Daniel Keysers, Skanda Koppula, Fangyu Liu, Adam Grycner, Alexey Gritsenko, Neil Houlsby, Manoj Kumar, Keran Rong, Julian Eisenschlos, Rishabh Kabra, Matthias Bauer, Matko Bošnjak, Xi Chen, Matthias Minderer, Paul Voigtlaender, Ioana Bica, Ivana Balazevic, Joan Puigcerver, Pinelopi Papalampidi, Olivier Henaff, Xi Xiong, Radu Soricut, Jeremiah Harmsen, and Xiaohua Zhai. 2024. PaliGemma: A versatile 3B VLM for transfer. arXiv preprint arXiv:2407.07726.
- Ali Furkan Biten, Ruben Tito, Andres Mafla, Lluis Gomez, Marçal Rusinol, Ernest Valveny, CV Jawahar, and Dimosthenis Karatzas. 2019. Scene text visual question answering. In ICCV.
- Victor Carbune, Hassan Mansoor, Fangyu Liu, Rahul Aralikatte, Gilles Baechler, Jindong Chen, and Abhanshu Sharma. 2024a. Chart-based reasoning: Transferring capabilities from llms to vlms. arXiv preprint arXiv:2403.12596.
- Victor Carbune, Hassan Mansoor, Fangyu Liu, Rahul Aralikatte, Gilles Baechler, Jindong Chen, and Abhanshu Sharma. 2024b. Chart-based reasoning: Transferring capabilities from llms to vlms. arXiv preprint arXiv:2403.12596.
- Paola Cascante-Bonilla, Hui Wu, Letao Wang, Rogerio S Feris, and Vicente Ordonez. 2022. Simvga: Exploring simulated environments for visual question answering. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5056-5066.
- Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. 2024. Seeclick: Harnessing gui grounding for advanced visual gui agents. Preprint, arXiv:2401.10935.
- Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini,

678

679

- 729
- 731 732

et al. 2024. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. arXiv preprint arXiv:2409.17146.

- Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. Scaling synthetic data creation with 1,000,000,000 personas. Preprint, arXiv:2406.20094.
- Boyu Gou, Ruohan Wang, Boyuan Zheng, Yanan Xie, Cheng Chang, Yiheng Shu, Huan Sun, and Yu Su. 2024. Navigating the digital world as humans do: Universal visual grounding for gui agents. arXiv preprint arXiv:2410.05243.
 - Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In Conference on Computer Vision and Pattern Recognition (CVPR).
 - Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. 2018. Vizwiz grand challenge: Answering visual questions from blind people. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 3608–3617.
- Yucheng Han, Chi Zhang, Xin Chen, Xu Yang, Zhibin Wang, Gang Yu, Bin Fu, and Hanwang Zhang. 2023. Chartllama: A multimodal llm for chart understanding and generation. arXiv preprint arXiv:2311.16483.
- Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. 2024. Cogagent: A visual language model for gui agents. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 14281–14290.
- Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 6700-6709.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. arXiv preprint arXiv:2310.06825.
- Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. 2017. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 2901-2910.
- Matthew Johnson-Roberson, Charles Barto, Rounak Mehta, Sharath Nittur Sridhar, Karl Rosaen, and Ram Vasudevan. 2017. Driving in the matrix: Can virtual worlds replace human-generated annotations for real world tasks? In 2017 IEEE International Conference

on Robotics and Automation (ICRA), page 746–753. **IEEE Press.**

733

734

735

736

738

739

740

741

742

743

744

745

746

747

748

749

753

754

755

756

757

758

759

760

761

762

764

765

766

767

768

769

771

772

774

775

776

777

778

779

780

781

782

783

784

785

- Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. 2018. DVQA: Understanding data visualizations via question answering. In CVPR.
- Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. 2017. FigureQA: An annotated figure dataset for visual reasoning. arXiv preprint arXiv:1710.07300.
- Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. 2016. A diagram is worth a dozen images. In ECCV.
- Yoonsik Kim, Moonbin Yim, and Ka Yeon Song. 2024. Tablevqa-bench: A visual question answering benchmark on multiple table domains. arXiv preprint arXiv:2404.19205.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024. LLaVA-OneVision: Easy visual task transfer. arXiv preprint arXiv:2408.03326.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In International conference on machine learning, pages 19730–19742. PMLR.
- Shengzhi Li and Nima Tajbakhsh. 2023. Scigraphqa: A large-scale synthetic multi-turn question-answering dataset for scientific graphs. arXiv preprint arXiv:2308.03349.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. Llavanext: Improved reasoning, ocr, and world knowledge.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. In NeurIPS.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In NeurIPS.
- Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and Ashwin Kalyan. 2023. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning. In ICLR.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. OK-VQA: A visual question answering benchmark requiring external knowledge. In CVPR.
- Ahmed Masry, Do Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In ACL.

787

- 837 838

- 841 842

- Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. 2022. InfographicVQA. In WACV.
- Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. DocVQA: A dataset for VQA on document images. In WACV.
- Nitesh Methani, Pritha Ganguly, Mitesh M Khapra, and Pratyush Kumar. 2020. PlotQA: Reasoning over scientific plots. In WACV.
- Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. 2019. Ocr-vqa: Visual question answering by reading text in images. In ICDAR.
- Srija Mukhopadhyay, Adnan Qidwai, Aparna Garimella, Pritika Ramu, Vivek Gupta, and Dan Roth. 2024. Unraveling the truth: Do VLMs really understand charts? a deep dive into consistency and robustness. In Findings of the Association for Computational Linguistics: EMNLP 2024, pages 16696–16717, Miami, Florida, USA. Association for Computational Linguistics.
- OpenAI. 2023. GPT-4 technical report. arXiv preprint arXiv:2303.08774.
- Ajay Patel, Colin Raffel, and Chris Callison-Burch. 2024. DataDreamer: A tool for synthetic data generation and reproducible LLM workflows. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3781–3799, Bangkok, Thailand. Association for Computational Linguistics.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In ICML.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of machine learning research, 21(140):1-67.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for SQuAD. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 784-789, Melbourne, Australia. Association for Computational Linguistics.
- N Reimers. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. arXiv preprint arXiv:1908.10084.
- Jonathan Roberts, Kai Han, Neil Houlsby, and Samuel Albanie. 2024. Scifibench: Benchmarking large multimodal models for scientific figure interpretation. arXiv preprint arXiv:2405.08807.

Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022. A-OKVQA: A benchmark for visual question answering using world knowledge. In ECCV.

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

885

887

888

889

890

891

892

893

894

- Risa Shinoda, Kuniaki Saito, Shohei Tanaka, Tosho Hirasawa, and Yoshitaka Ushiku. 2024. Sbs figures: Pre-training figure qa from stage-by-stage synthesized images. arXiv preprint arXiv:2412.17606.
- Noah Siegel, Zachary Horvitz, Roie Levin, Santosh Divvala, and Ali Farhadi. 2016. Figureseer: Parsing result-figures in research papers. In Computer Vision-ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII 14, pages 664-680. Springer.
- Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus Rohrbach. 2019. Towards VQA models that can read. In CVPR.
- 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.
- Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. 2024. To cot or not to cot? chain-of-thought helps mainly on math and symbolic reasoning. arXiv preprint arXiv:2409.12183.
- Gemini Team. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530.
- Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. 2024. Cambrian-1: A fully open, vision-centric exploration of multimodal LLMs. In NeurIPS.
- Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. 2021. Multimodal few-shot learning with frozen language models. Advances in Neural Information Processing Systems, 34:200–212.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837.
- Renqiu Xia, Bo Zhang, Hancheng Ye, Xiangchao Yan, Qi Liu, Hongbin Zhou, Zijun Chen, Min Dou, Botian Shi, Junchi Yan, et al. 2024. Chartx & chartvlm: A versatile benchmark and foundation model for complicated chart reasoning. arXiv preprint arXiv:2402.12185.

Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan
Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu,
Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. 2024. Osworld: Benchmarking multimodal agents for openended tasks in real computer environments. *Preprint*, arXiv:2404.07972.

904

905

906 907

908

909

910

911

912 913

914

915

916

917 918

919

921

- Zhengzhuo Xu, Sinan Du, Yiyan Qi, Chengjin Xu, Chun Yuan, and Jian Guo. 2023. Chartbench: A benchmark for complex visual reasoning in charts. *arXiv preprint arXiv:2312.15915*.
- Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu, Yutong Dai, Michael S Ryoo, et al. 2024. xGen-MM (BLIP-3): A family of open large multimodal models. *arXiv preprint arXiv:2408.08872*.
 - Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng, Alexander Ratner, Ranjay Krishna, Jiaming Shen, and Chao Zhang. 2023. Large language model as attributed training data generator: A tale of diversity and bias. *Preprint*, arXiv:2306.15895.
- Jieyu Zhang, Le Xue, Linxin Song, Jun Wang, Weikai Huang, Manli Shu, An Yan, Zixian Ma, Juan Carlos Niebles, Caiming Xiong, et al. 2024. Provision: Programmatically scaling vision-centric instruction data for multimodal language models. *arXiv preprint arXiv:2412.07012*.

Implementation Details Α

A.1 Prompts

We provide the prompt templates in Figure 8 for the HTMLDocumentPipeline as an example to illustrate the prompts used across our code-guided synthetic data generation pipelines.

Topic Generation: You are an expert in document generation and have a broad knowledge of different topics.

My persona is: "PERSONA" I want you to generate NUM_TOPICS topics for FIGURE_TYPE that I will be interested in or I may see during my daily life given my persona.

Here are the requirements Each topic is a high-level summary of the contents in FIG-URE_TYPE with some design details, e.g., "the utility bill for the month of January 2022 with a detailed breakdown of charges" 2. The topics should be diverse to help me generate varied documents. Each topic should be unique and not overlap with others.

3. The topics are conditioned on the document type. Please ensure the topics you provided can be best visualized in "FIGURE TYPE" 4. All topics must be in English, even if the persona is non-English. 5. List NUM_TOPICS topics for "PERSONA" and separate them with a | character, e.g., topic1 | topic2 | | topicN.

Do not include any additional text at the beginning or end of your response.

Data Generation: You are an expert in content creation and have broad knowledge about various topics.

My persona is: "PERSONA" I need some materials about "TOPIC", which can be used to generate a FIGURE_TYPE. Here are the requirements:

1. The materials should be related to the topic and customized according to my persona. Its structure must be suitable for the FIGURE_TYPE

2. The materials should be realistic, and the contents should be named using real-world entities. Do not use placeholder names like xxA, xxB, etc. Do not use template data like [Name], [Date], etc. 3. The materials should be diverse and contain information from different aspects of the topic to ensure the document is informative.

4. Do not provide too many materials. Just provide key pieces of information that are essential for a **one-page document 5. All materials must be in English, even if the persona is non-

English. Please provide the materials in JSON format without additional text

at the beginning or end.

Code Generation: You are an expert web designer and are good at writing HTML to create documents. My persona is: "PERSONA" I have some materials about TOPIC

which can be used to generate a FIGURE_TYPE. Here are the materials (JSON format):

<data> DATA </data>

Please use HTML and CSS to generate a FIGURE_TYPE using the data provided.

Here are the requirements: 1. **Style Requirements**: Feel free to use any CSS framework, libraries, JavaScript plugins, or other tools to create the document. (1) Try to be creative and make the web page style, fonts, colors, borders and visual layout unique with CSS. Taking persona, topic, and document type into consideration when designing the docu ment.

(2) Select the appropriate design scale (e.g., margins, page size, layout, etc) to ensure the information in the document is clear and easy to understand, with no text overlapping, etc.

Do not make the page too long or too sparse.** All contents should be in **one page**. This is very important. 2. **Code Requirements**:

(1) You need to hardcode the provided data into the HTML script to generate the document. Be careful with the syntax and formatting of the HTML

(2) Put everything in one HTML file. Do not use external CSS or JavaScript files.

3. **Output Requirements**: Put "thtml at the beginning and "t at the end of the script to separate the code from the text.

Please don't answer with any additional text in the script, your whole response should be the HTML code which can be directly executed.

Instruction Generation: You are an expert in data analysis and good at asking questions about documents. My persona is: "persona" I want you to generate some guestion-answer pairs of a FIGURE_TYPE about TOPIC, which I would ask. Instead of showing the document, I provide the data and the code that generates the document. <data> DATA </data> <code> CODE </code> Please come up with a list of *reasonable questions* that people will ask when they see the rendered document. Here are the requirements 1. **Question Types**: All questions should be short-answer questions that are answerable based on the visual information in the document. All questions can be answered with a single word, phrase, or number. (as short as possible) (1) **Information Retrieval questions** ask for specific information in the document, such as numbers, names, dates, titles, etc. The questions should cover different aspects (areas) of the document. This is the most common type of question. (2) **Reasoning questions** require reasoning over multiple information in the document. These questions should be more challenging and require a deeper understanding of the document (3) **Document Type-specific questions** are questions that are specific and unique to this document type FIGURE_TYPE. These questions should be tailored to the content and structure of the document. 2. **Response Format**: Use | to separate the question, explanation, and concise answer for each example. (1) Follow this format: question | explanation | concise answer, e.g., what is the total revenue? | The total revenue is the sum of all revenue sources in the document, which is \$2000 + \$3000 + \$5000 \$10000. | \$10000 (2) Separate the question-answer pairs by double newlines. question1 | explanation1 | answer1 question2 | explanation2 | answer2... (3) Do not provide too many questions, 5-10 questions are enough. Focus on the diversity and quality of the questions. Try to cover different aspects of the document. (4) The concise answer should be as short as possible and directly answer the question. The answer should be faithful and exactly the same as what you would expect to see in the document, don't rephrase it. All words in the answer should be processed in natural language, no coding terms/characters

Please follow the format strictly and do not include any additional text at the beginning or end of your response.

Figure 8: Prompt templates used for HTML Document Pipeline, including all four stages of generation: topic, data, code, and instruction.

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

Rendering Tools and Pipelines A.2

We design 20 generation pipelines built on 11 rendering tools to support the creation of nine categories of text-rich images:(1) Charts: Matplotlib VegaLite, Plotly, LaTeX, HTML; (2) Documents: LaTeX, HTML; (3) Tables: LaTeX, Matplotlib, Plotly, HTML; (4) Diagrams: Graphviz, LaTeX, Mermaid; (5) Math Problems: LaTeX; (6) Vector Graphics: SVG, Asymptote; (7) Music Sheets: LilyPond; (8) Electrical Circuits: LaTeX; (9) Chemical Structures: Rdkit. In addition, we implement a separate pipeline for generating pointing data using HTML as the rendering tool.

Queries to Construct CoSyn-400K A.3

Since CoSyn accepts textual queries to control data generation, we use a diverse set of queries for each type of text-rich image to ensure broad domain coverage. Below are some examples of the queries used to generate CoSyn-400K:

929

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

• Charts: Bar, Line, Pie, Diverge bar, Bubble, Scatter, Histogram, Area, Box plot, Heatmap, Error bar, Radar chart, Rose chart, Stem plot, Stairs plot, Violin chart, 2D contour, Distplots, Log plot, Ternary plots/contour, Candlestick charts, Time series, etc. (51 queries in total)

950

951

953

954

955

957

959

960

961

962

963

965

967

969

972

973

974

975

977

978

979

983

984

985

987

992

996

997

1000

- **Documents**: Letter, Form, Report, Receipt, Invoice, Restaurant menu, Newsletter, Schedule, Manual, Brochure, Transaction document, Agenda, Memo, Financial report, Telephone records, Note, Budget, Meeting minutes, Bill, Catalog, Email, Fax, Policy document, Resume, Infographics, Process infographic, Statistical infographic, etc. (**107** queries in total)
 - Math Problems: Algebra, Counting, Probability, Geometry, Number theory, Precalculus, Prealgebra, Intermediate Algebra, Statistics, Functions, Complex numbers, Logarithms, Inequalities, Linear equations, Exponents, Series, College Algebra, Calculus, Advanced calculus, Linear algebra, Solid geometry, Analytic geometry, Polynomial arithmetic, etc. (110 queries in total)
 - **Tables**: Financial table, Simple table, Pivot table, Comparison table, Timeline table, Decision table, Truth table, Lookup table, Periodic table, Statistical table, Timetable, Hierarchical table, Matrix table, Contingency table, Logarithmic table, Correlation table, etc. (**35** queries in total)
 - **Diagrams**: Flow chart, Directed graph, Undirected graph, Decision tree, Mind map, Gantt charts, Finite state machine, Quadrant chart, Chord diagrams, Network diagrams, Sankey diagram, Entity relationship diagram, Sequence diagrams, Bottom-up flow chart, Timeline, State diagram, Concept map, Family tree, Programming flowchart, etc. (**34** queries in total)
 - Vector Graphics: Visual intelligence test, Spatial intelligence test, Geometry, Solid geometry, Analytic geometry, Polynomial graphs, Trigonometry, Polar coordinates, Coordinate system, Topology, Graph theory, Plane geometry, Functions, Calculus, Vectors, Angles, Perimeter and area problems, etc. (36 queries in total)
 - Sheet Music: Classical, Pop, Rock, Jazz, Blues, Hip Hop, Rap, Electronic, Country, Folk, Rhythm and blues, Soul, Reggae, Metal, Punk, Theme, Dance, etc. (34 queries in total)
 - Electrical Circuits: Series, Parallel, Hybrid, Household appliances, Industrial appliances, Mobile device, Low-power appliances, High-power appliances, etc. (**30** queries in total)
 - Chemical Structures: Drug, Organic, Inor-

ganic, Protein, Acids, Bases, Gases, Liquids,1001Solids, Oxidizers, Flammable liquids, Toxic1002chemicals, Hazardous chemicals, Aromatic compounds, Aliphatic compounds, Polymers, Metals,1004Alloys, Electrolytes, etc. (100 queries in total)1005

A.4 Academic Datasets

During the supervised fine-tuning stage, we include academic datasets in addition to our synthetic datasets. Below, we provide details on the size of these datasets and the evaluation metrics used.

Dataset Size. The number in parentheses indicates the number of training images for each dataset: ChartQA (28.3K), DocVQA (39.5K), InfographicVQA (23.9K), AI2 Diagrams (11.4K), TextVQA (34.6K), VQAv2 (82.8K), GQA (72.1K), OK-VQA (9.0K), OCR-VQA (166.0K), A-OKVQA (17.1K), ScienceQA (6.2K), TabMWP (23.1K), ST-VQA (18.9K), TallyQA (133.0K), DVQA (200.0K), FigureQA (100.0K), PlotQA (160.0K). We downsample some very large synthetic datasets, such as DVQA, FigureQA, and PlotQA, to balance the dataset size. In total, we use approximately 1.1M images from academic datasets.

Evaluation Metrics. We adopt their official evaluation metrics for the seven text-rich datasets. For ChartQA, we use relaxed correctness, which allows a 5% difference for float number answers. For DocQA and InfoQA, we report Average Normalized Levenshtein Similarity (ANLS). For TableVQA, we report the average performance across the four subsets (VTabFact, VWTQ, VWTQ-Syn, FinTabNetQA) using the metrics provided in this repo. We report the multiple choice accuracy for AI2D, VQA score (Goyal et al., 2017) for TextVQA, and SQuAD F1 score (Rajpurkar et al., 2018) for ScreenQA.

A.5 Training Details

Image Preprocessing. We adopt the same image preprocessing as Molmo (Deitke et al., 2024), where each input image is cropped into multiple overlapping crops before being encoded by CLIP. During training, we limit the maximum number of crops to 12, but we increase it to 25 at testing time to accommodate the high resolution of text-rich images. This strategy boosts the inference performance without increasing training costs.

Hyper Parameters. We set the maximum sequence length for training is 2304 tokens. We use the same learning rate of 1e-6 for the MLP con-

nector, LLM, and visual encoder, with batch size 1050 32. The best-performing model is trained for 60K 1051 steps with 200 warm-up steps and a cosine sched-1052 uler with an end factor of 0.1. All experiments are 1053 run on a single TPU v3-128. 1054

B **Additional Analysis**

1055

1057 1058

1059

1060

1061

1062

1064

1067

1068

1069

1070

1071

1073

1074

1075

1076

1078

1079

1080

1081

1082

1083

1084

1086

I

We conduct additional analyses below to investigate further why our synthetic data can effectively enhance vision-language models.

Our synthetic data is more diverse. To quantify the diversity of images and text in our synthetic dataset $\mathcal{D} = \{(I, T)\}$, we propose the following two metrics to compute the diversity:

$$\mathbf{Diversity}(\mathcal{D})_{\mathbf{Image}} = \frac{1}{|\mathcal{D}|^2 - |\mathcal{D}|} \sum_{I_i \in \mathcal{D}} \sum_{I_j \in \mathcal{D}}^{i \neq j} \left(1 - \sin(I_i, I_j)\right) \quad (1)$$

 $\text{Diversity}(\mathcal{D})_{\text{Text}} = \frac{1}{|\mathcal{D}|^2 - |\mathcal{D}|} \sum_{T_i \in \mathcal{D}} \sum_{T_i \in \mathcal{D}}^{i \neq j} \left(1 - \sin(T_i, T_j)\right) \quad (2)$

where $sim(\cdot)$ is the cosine similarity function. Both metrics compute the average pairwise cosine distance between the features of every instance in the dataset. For image diversity, we extract features using CLIP, while for text diversity, we use Sentence-BERT (Reimers, 2019) to obtain embeddings of question-answer pairs. Table 4 shows that our synthetic charts are significantly more diverse than those in existing datasets, such as FigureQA and ChartQA, in both image and text diversity.

Dataset	Image Diversity	Text Diversity
FigureQA	0.268	0.567
DVQA	0.307	0.752
PlotQA	0.420	0.743
ChartQA	0.340	0.742
Ours (Charts)	0.596	0.823

Table 4: Compare image and text diversity across different chart datasets. We randomly sample 10K instances from each dataset to compute the results.

Diversity correlates with model performance. We observe that data diversity significantly affects model performance on downstream tasks. To investigate this, we compare synthetic chart data generated using only a single tool (Matplotlib) with charts generated by all five tools available in our CoSyn system. As shown in Table 5, using multiple tools results in higher image diversity and notably improved performance on ChartQA. This experiment underscores the importance of data diversity for enhancing the generalizability of models.

n. of Tools	Diversity		ChartQA	
		Average	Machine	Human
Single Multiple	0.572 0.607	73.9 75.2	66.5 68.6	81.5 82.0

Table 5: Single vs. Multiple Rendering Tools for Data **Generation.** Each row uses the same number of 45K synthetic images. Single only uses Matplotlib, while Multiple involves four other rendering tools: HTML, LaTex, Plotly, and VegaLite.



Figure 9: Scaling the Size of Synthetic Data. We evaluate the zero-shot performance on ChartQA of models fine-tuned on increasing numbers of synthetic images.

Scaling the size of synthetic data. In addition to diversity, the scale of synthetic data also impacts model performance. As shown in Figure 9, increasing the number of synthetic chart images leads to improved performance on ChartQA. This demonstrates that scaling up synthetic data can further enhance VLMs on downstream tasks. Due to resource constraints, our final dataset consists of 400K images, which cost us about \$8,000. Future work could explore scaling up the dataset size to push the boundaries of synthetic data's potential.

LLM for Data Generation		ChartQA	
	Average	Machine	Human
GPT-40	72.4	65.8	78.9
Claude-3.5-sonnet	77.2	71.0	83.8

Table 8: Compare the LLMs used for synthetic data generation. For both LLMs, we create 100K synthetic charts for fine-tuning the VLMs. We report the zeroshot evaluation results on ChartQA.

Compare LLMs for synthetic data generation. 1098 In the default setting, CoSyn uses Claude-3.5-1099 sonnet as the underlying LLM for code generation. 1100 To highlight the importance of strong coding capabilities, we compare it with data generated by 1102 GPT-40. As shown in Table 8, synthetic data gen-1103 erated by Claude-3.5-sonnet yields significantly 1104

1101

Prompt Type	ChartQA	DocVQA	InfoVQA	TableVQA	AI2D	TextVQA	ScreenQA	NutritionQA
CoT	86.3	87.4	63.8	65.8	86.0	70.9	79.0	76.0
Short Answer	83.1	90.0	70.5	64.3	91.9	82.0	80.1	62.0

Table 6: Alation of using chain-of-thought (CoT) in prompts. CoT means letting the model provide reasoning steps before giving the final answer. Short Answer prompts the model to answer with as few words as possible.

FT Data	ChartQA	DocVQA	InfoVQA	$\textbf{TableVQA}^{\dagger}$	AI2D	TextVQA	$\textbf{ScreenQA}^{\dagger}$	Average
Aux only*	60.7	56.2	39.7	43.1	81.7	68.5	61.3	58.7
Syn only*	79.4	80.5	60.1	64.4	68.6	63.6	76.6	70.5
Aux + Syn*	80.8	82.9	59.8	64.9	83.9	72.7	78.1	74.7
Eval only	77.4	87.4	63.8	51.8	91.3	81.1	78.1	75.9
Eval + Aux	81.4	87.9	68.2	53.6	91.6	81.8	77.0	77.3
Eval + Aux + Syn	86.3	90.0	70.5	65.8	91.9	82.0	80.1	80.9

Table 7: Alation of the data selection for supervised fine-tuning. Aux, Syn, and Eval stand for auxiliary, synthetic, and evaluation datasets, respectively. The rows with * represent zero-shot models (without using any training examples from any of the evaluation datasets). The datasets with [†] are test-only datasets (no training splits), which means all numbers on these datasets are zero-shot performance.

better results than GPT-40. Our qualitative observation reveals that GPT-40 has a higher failure rate in code generation, particularly for less common coding languages or libraries. This result emphasizes that a strong LLM is essential for the successful synthetic data generation for VLMs.

1105

1106

1107

1108

1109

1110

1111Quantify the contributions of synthetic data. Ta-1112ble 7 presents the performance across benchmarks1113using different combinations of supervised fine-1114tuning data. A clear trend shows that synthetic1115data significantly contributes in both zero-shot and1116supervised settings. Adding our synthetic data con-1117sistently boosts performance on each benchmark.

The impact of Chain-of-thought reasoning. We 1118 compare the performance of CoT and short-answer 1119 prompts in Table 6. CoT reasoning improves perfor-1120 mance on ChartQA, TableVQA, and NutritionQA, 1121 1122 where questions require multi-hop and mathematical reasoning that aligns with the findings in lan-1123 guage tasks (Sprague et al., 2024). However, short-1124 answer prompts yield better results for the other 1125 five datasets due to their annotation biases favoring 1126 concise responses. CoT responses tend to be more 1127 verbose, which may not match the ground-truth 1128 answers exactly, resulting in a performance drop. 1129

Document Pointing Task. To further validate 1130 the effectiveness of our synthetic pointing data, 1131 we introduce DocPointQA, a new pointing task 1132 1133 with 300 question-point pairs annotated from the DocVQA validation set (Figure 11). We compare 1134 models trained on human-annotated PixMo-point 1135 data (155K examples), our synthetic pointing data 1136 (65K examples), and their combination. Since 1137

DocPointQA requires multiple-point answers, we report precision, recall, F1 score, and L2 distance (lower is better) after mapping predicted points to ground truth, following the same setup as Molmo (Deitke et al., 2024). As shown in Table 9, the model trained on our synthetic data outperforms the one trained on PixMo-point. Performance improves even further when both datasets are combined, demonstrating the effectiveness of synthetic data in enhancing the pointing capabilities of vision-language models.

Pointing Data	Precision	Recall	F1	Distance \downarrow
PixMo-point	49.7	49.3	52.7	17.3
Synthetic (Ours) Combined (Ours)	63.8 69.9	66.1 70.6	62.8 70.7	9.2 8.8

Table 9: **Zero-shot Pointing on DocPointQA.** We compare the models trained on different pointing data. Combined stands for combining PixMo-point (human-annotated) (Deitke et al., 2024) with our synthetic data.

Figure 10 and 11 show the examples from our an-

notated NutritionQA and DocPointQA. Figures 12

- 18 list examples from the 9 categories of synthetic

text-rich images. Figure 19 illustrates examples

Use Of AI Assistants. We use AI to fix some typos

C Qualitative Examples

from the synthetic pointing dataset.

and grammar. Authors write all contents.

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1150

1151 1152

1153

1154



Figure 10: Examples from our newly collected NutritionQA dataset.



Figure 11: Examples from our newly collected **DocPointQA** dataset.



Figure 12: Randomly selected examples from our synthetic chart data.



Figure 13: Randomly selected examples from our synthetic document data.

Category		Ye	ear	
	2020	2021	2022	2023
Engine Components	1500	1800	2000	2200
Body Parts	1200	1400	1600	1700
Interior Trim	900	950	1200	1300
: Which year ha	d the l	lowes	t sales	for
nterior Trim? A:	2020			
• Which month	had th	e hiøł	nest cu	stomer
ot traffic? A: De	ecemb	er.	1050 00	stomer

Figure 14: Randomly selected examples from our synthetic table data.



Figure 15: Randomly selected examples from our synthetic math data.



Figure 16: Randomly selected examples from our synthetic diagram data.



Figure 17: Randomly selected examples from our synthetic vector graphic data.



Figure 18: Randomly selected examples from our synthetic sheet music, circuits and chemical structures.



Figure 19: Randomly selected examples from our synthetic pointing data.