

MetaSynth: Meta-Prompting-Driven Agentic Scaffolds for Diverse Synthetic Data Generation

Anonymous ACL submission

Abstract

Recent smaller language models such as Phi-3.5 and Phi-4 rely on synthetic data generated using larger Language models. Questions remain about leveraging synthetic data for other use cases, such as adapting LLMs to specific domains. A key limitation of synthetic data is *low diversity*, which negatively impacts its downstream applicability for improving other models. To address this, we propose METASYNTH, a method for generating synthetic data that enhances diversity through meta-prompting, where a language model orchestrates multiple “expert” LLM *agents* to collaboratively generate data. Using only **25 million** tokens of synthetic data generated with METASYNTH, we successfully adapt a well-trained LLM (Mistral-7B) to two specialized domains—Finance and Biomedicine—without compromising the capabilities of the resulting model in general tasks. In addition, we evaluate the diversity of our synthetic data using seven automated metrics, and find that it approaches the diversity of LLM pre-training corpora.

Continually pre-training Mistral-7B with MetaSynth notably outperforms the base LLM, showing improvements of up to 4.08% in Finance and 13.75% in Biomedicine. The same model shows degraded performance when trained on data generated using a template prompt, even when the template includes prior generations and varying In-Context exemplars of real data. Our findings suggest that a few million tokens of diverse synthetic data without mixing any real data, is sufficient for effective domain adaptation when using MetaSynth.

1 Introduction

Human generated public text data cannot sustain the continued scaling and expansion of large language models (LLMs). It has been

argued by Villalobos et al. (2024) that the available stock of public human text data will be fully utilized by 2028 if current LLM development trends continue, or earlier if LLMs are trained on more data than is compute optimal. This is evidenced in e.g., Llama 3 (Grattafiori et al., 2024) which uses one order of magnitude more data compared to only two year old estimates of compute optimal large language models (Hoffmann et al., 2022). As a potential remedy, synthetic data generated with LLMs has shown remarkable potential to alleviate the impending issue of data scarcity for future model scaling.

However, low diversity is a key issue in any type of synthetic data. In this work, we hypothesize that there are two prominent reasons that affect the diversity of data synthesized by LLMs: a) the choice of seed instances used to initialize data generation and b) the prompts used, which commonly follow predefined templates, where variation in the prompt is mainly introduced via placeholders whose content is populated dynamically. Examples of data generation methods which use template-like prompts with variation include: *Self-prompting* (Li et al., 2024), *Attrprompt* (Yu et al., 2023), *CLINGEN* (Xu et al., 2025) and *Explore-Instruct* (Wan et al., 2023), among others. We contend that this variation yields limited diversity. For instance, generating a collection of domain-specific (e.g., financial) texts with similar prompts results in repetitive sentence structures—many texts begin with lexical patterns such as “In today’s ever-changing financial landscape” or “as the financial world evolves”—and often contain recurring phrases, and generic buzzwords.

Recently, Suzgun and Kalai (2024) find that *Meta-prompting* (Zhang et al., 2024) approaches – where an LLM itself writes the

prompts to solve a problem – can elicit more diverse and creative outputs, significantly improving problem-solving capabilities for mathematical and algorithmic reasoning tasks, largely due to the feedback, self-verification, chain-of-thought (Wei et al., 2023), and planning dynamics inherent in these approaches. It has also been shown that using an optimized meta-prompt can improve the quality and downstream effectiveness of LLM generated synthetic data (Kim et al., 2024). We argue that a key use case for synthetic data arises when abundant real data exists in the form of pre-training corpora, but one wishes to effectively tailor an LLM to a specific domain using only a small amount of carefully generated synthetic data. In this work, we investigate *data-efficient domain adaptation* through meta-prompting, where a language model is instructed to act as a supervisor that writes specialized prompts for other models to collaboratively generate diverse data. Our contributions are as follows:

- (1) We propose METASYNTH, a method to create diverse synthetic documents for continual pre-training (CPT) by leveraging a meta language model (which we refer to as meta-LM) and *Conditional Instance Generation* – where the meta-LM categorizes, and keeps track of each generated instance in memory, to ensure distinctness between them.
- (2) We propose METASYNTH-INSTRUCT, which can generate and iteratively evolve complex instructions for *instruction pre-training*. Notably, this evolution is entirely driven by prompts written by the meta-LM itself. Furthermore, unlike other instruction-pretraining approaches e.g., ((Cheng et al., 2024a), (Cheng et al., 2024b)) our instructions are purely evolved from contexts synthesized by METASYNTH i.e., without using any human-written text (section 3.3).
- (3) METASYNTH-INSTRUCT can also synthesize training data for fine-tuning encoder models such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019) etc. We observe that encoders fine-tuned on this data can outperform those fine-tuned on data generated with template-based prompting (section 6).
- (4) We generate synthetic documents for continual pre-training by prompting an LLM with its prior outputs (memory) where the prompt follows a predefined template containing in-

context exemplars of real data. However, when these synthetic documents are mixed with real data in a 1:1 ratio over 25M tokens, we find that it *does not* improve the Mistral-7B base model, and even leads to slight performance degradation across two domains. In contrast, using 25M tokens of diverse synthetic data from MetaSynth yields substantial improvements to the base model across various ratios of mixing real and synthetic data. Experiments on ten datasets in Finance and Biomedicine—evaluating nine mixing ratios following Cheng et al. (2024b)—indicate that **mixing real data with synthetic data is not needed if synthetic data is diverse**. (section 5).

(6) We systematically measure the diversity of LLM generated synthetic data across multiple dimensions using seven automated metrics, including the *Task2Vec* diversity coefficient (Lee et al., 2023), which encapsulates formal notions of data diversity, among others (see Section 4). We find that our approach significantly improves the diversity of generated data relative to template-based prompting (section 4). We argue that model degradation and “model collapse” (as discussed, inter alia, in (Shumailov et al., 2024; Seddik et al., 2024; Gerstgrasser et al., 2024)) can be avoided even when training solely on a small amount of synthetic data if it is sufficiently diverse. We present our method below, which we view through the lens of inference-time compute scaling, to ensure diversity in synthetic data.

2 Meta-Prompting

2.1 Meta-LM

At a high level our procedure for synthesizing a diverse collection of documents leverages two ideas: Meta-Prompting (Suzgun and Kalai, 2024; Zhang et al., 2024) and *Conditional Instance Generation* (refer to section 2.1). Meta-prompting leverages a central meta-LM to coordinate and execute multiple independent inquiries and subsequently synthesize their responses to render a final response. This is realized via a high-level “meta” prompt which instructs an LM to break down complex tasks (such as generating a diverse collection of documents) into smaller or more manageable subtasks. Each of these subtasks is assigned to spe-

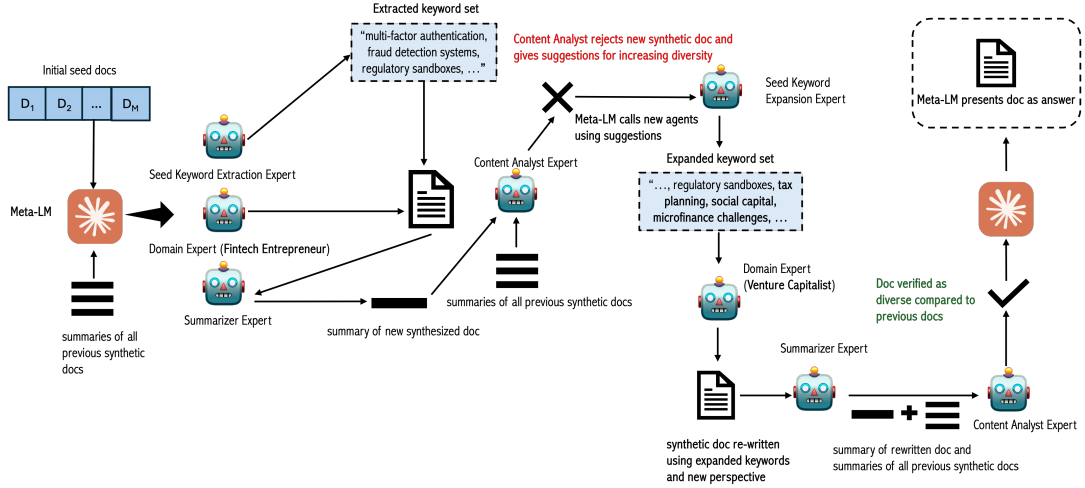


Figure 1: Demonstration of an example METASYNTH agentic workflow for synthesizing a financial document. A meta-LM orchestrates various expert agents that iteratively refine and generate diverse documents conditioned on an initial set of seed documents and previously synthesized documents. Refer to Section 2.2 for a detailed description of the workflow.

cialized expert models (also known as “agents”) where both the choice of the agent and the instructions to the agent are dynamically composed by the meta-LM depending upon the nature of the subtask. In this work, we adapt the task agnostic meta-prompt from Suzgun and Kalai (2024) to specifically focus on generating diverse synthetic data. The meta-LM serves as an orchestrator overseeing communication between these agents in a centralized multi-agent system (MAS) (Guo et al., 2024), where the agents cannot directly interact with each other; and also carries forward the thread of the process by applying its own critical thinking, reasoning and verification skills throughout. Further, to enable conditional instance generation (section 2.1), the meta-LM is equipped with memory to become stateful—a message history comprising its own responses (which include the selection of agents and formulation of instructions for them) and the responses from various agents. *Only the meta-LM* has access to the complete history, while the agents it invokes are limited to selectively shared information, seeing only what the meta-LM chooses to share with them. Being provided with only partial information pertaining to the task to solve, allows an agent to consider new perspectives with “fresh eyes” (Suzgun and Kalai, 2024) and potentially correct the meta-LM’s errors. In this work, we use Claude 3 Sonnet (Anthropic, 2024) as the

meta-LM.

We further motivate the need for agentic scaffolding by drawing an analogy to multi-disciplinary problem solving: complex tasks are often best addressed by leveraging diverse expertise rather than relying on a single, monolithic approach. As shown by Wu et al. (2023); Yao et al. (2023), distinct agents can specialize in decision making, problem decomposition, and mitigating issues such as error propagation in chain-of-thought reasoning. To generate a single synthetic instance (e.g., document or instruction), the meta-LM can invoke agents arbitrarily. However, to ensure that: a) each synthetic instance (document or instruction) is sufficiently distinct from all previously generated instances, and b) the meta procedure does not degenerate: we impose specific constraints within the meta-prompt which specify that certain types of agents must always be invoked, accompanied by an in-context exemplar that demonstrates the invocation process for those agents. The required agents depend on the task—whether synthesizing documents, instructions, or instances from an existing dataset (refer to the meta-prompts in Appendix J and K). Without these constraints,¹ the procedure risks degenerative loops, where repetitive exchanges between a meta-LM and

¹Even with these constraints, degeneration can still occur due to noisy message passing between the meta-LM and agents.

agent(s) may hinder task completion.

Beyond this, the procedure remains open-ended, allowing the meta-LM to leverage any type of agent to enhance instance diversity in the final set of synthesized instances. In this work, the meta-LM (Claude 3 Sonnet) also serves as the agent LM, though any advanced instruction-following models can fulfill these roles.

Conditional Instance Generation Our document synthesis approach relies on a continuously expanding *instance classification table* appended to the meta-LM’s history in each iteration, tracking and categorizing generated instances. New instances (documents or instructions) are conditioned on prior ones to ensure distinctness while adhering to seed keywords or documents. This process is guided by two agents: the “Seed Keyword Expansion Expert” and the “Content Analyst Expert.” Documents are compared via summaries (generated by the “Summarizer Expert”), while instructions are compared directly. Summarization mitigates LLM context window limitations when managing a large set of prior documents. The Content Analyst Expert suggests diversity-enhancing modifications, such as expanding the seed keyword set with related terms or incorporating new personas. Conditional instance generation is illustrated in algorithm 2, appendix B, and the meta-prompting procedure that we adapt from Suzgun and Kalai (2024) is shown in algorithm 1, appendix B.

Exit Criteria and Error Handling At each iteration, the meta-LM, conditioned on its history, must either call an agent or return a final response marked by the <end> token (indicating that the desired number of instances have been synthesized from the initial seeds). Otherwise, an error is appended to its history and the model is prompted to retry. After N attempts the iteration is discarded.

2.2 Execution

Figure 1 illustrates “fresh eyes” and “conditional instance generation” through an example execution history that synthesizes a new financial domain document given initial seed documents and previously synthesized documents: (1) The meta-LM consults a “Seed Keyword Extraction Expert,” a “Domain Expert,” and a “Summarizer Expert.” (2) The Seed Keyword

Extraction Expert extracts representative keywords (e.g., “multi-factor authentication,” “fraud detection,” “regulatory sandboxes”), which the meta-LM uses to instruct a Domain Expert (e.g., a “Fintech Entrepreneur”) to generate a document. (3) The Domain Expert writes the document, which the Summarizer Expert condenses before the meta-LM accepts it. (4) The meta-LM then instructs the Domain Expert to generate a second document that adheres to the same keywords while differing in content and style. (5) To verify diversity, the meta-LM consults the Summarizer Expert and a “Content Analyst Expert.” If the second document is deemed too similar to the first, the Content Analyst provides feedback. (6) In response, the meta-LM calls a “Seed Keyword Expansion Expert” to enrich the keyword set and instructs a new Domain Expert (e.g., a “Venture Capitalist”) to rewrite the document from a fresh perspective. (7) The Summarizer and Content Analyst Experts reassess the revised document, and a “Writing/Linguistic Expert” may be consulted for stylistic diversity. Once confirmed as sufficiently distinct, the document is accepted, and the process continues for generating subsequent documents.

3 Synthetic Data Generation

3.1 Baseline: Template Prompting

We introduce a strong baseline for synthetic data generation that uses a static template-based prompt (refer to appendix I) with a placeholder populated by five-shot examples of real documents randomly selected from a domain specific subset of Common Crawl². Additionally, this generation process is also *conditional*, as the data generator is equipped with memory, allowing it to reference previously generated documents while being instructed to ensure that each new document remains distinct from prior outputs. Refer to appendix I for examples of documents synthesized with template prompting.

3.2 Meta-Synth: Synthetic Document Generation

Random Seed Selection For generating synthetic documents, we propose two methods for selecting a set of seed instances. The first is

²<https://commoncrawl.org/>

Setting	Compression Ratio ↓	Task2Vec Div. Coeff ↑	Remote Clique ↑	Chamfer Distance ↑	1-GD ↑	4-GD ↑	MIF ↑
Template Prompting	<u>3.6674</u>	<u>0.1576</u>	<u>0.1964</u>	<u>0.0897</u>	<u>0.0198</u>	<u>0.9224</u>	<u>8.5614</u>
Common Crawl	2.7380 (-25.34%)	0.212 (+34.52%)	0.3036 (+54.58%)	0.2359 (+162.99%)	0.0621 (+213.64%)	1.6080 (+74.33%)	8.1263 (-5.08%)
Synth. Data (Seed Keywords)	3.4443 (-6.08%)	0.1757 (+11.49%)	0.2191 (+11.56%)	0.1351 (+50.61%)	0.0345 (+74.24%)	1.1749 (+27.37%)	9.0016 (+5.14%)
Synth. Data (Seed Documents)	3.1495 (-14.12%)	0.1788 (+13.45%)	0.2047 (+4.23%)	0.1383 (+54.18%)	0.0390 (+96.97%)	1.3468 (+46.01%)	8.9150 (+4.13%)
Wikipedia	2.6088 (-24.82%)	0.1892 (+20.05%)	0.2868 (+46.03%)	0.2416 (+169.34%)	0.1046 (+428.28%)	1.6997 (+84.27%)	8.3149 (-2.88%)

(a) Evaluating diversity metrics of synthetic data generation methods from **finance** domain.

Setting	Compression Ratio ↓	Task2Vec Div. Coeff ↑	Remote Clique ↑	Chamfer Distance ↑	1-GD ↑	4-GD ↑	MIF ↑
Template Prompting	<u>3.4699</u>	<u>0.1575</u>	<u>0.2295</u>	<u>0.1056</u>	<u>0.0278</u>	<u>1.0035</u>	<u>8.7463</u>
Common Crawl	2.6717 (-23.00%)	0.2068 (+31.30%)	0.3130 (+36.38%)	0.2451 (+132.10%)	0.0703 (+152.88%)	1.6524 (+64.66%)	8.2744 (-5.40%)
Synth. Docs (Seed Keywords)	3.1537 (-9.11%)	0.1760 (+11.75%)	0.2277 (-0.79%)	0.1426 (+35.04%)	0.0403 (+44.96%)	1.3323 (+32.77%)	8.9503 (+2.33%)
Synth. Docs (Seed Docs)	3.0649 (-11.67%)	0.1793 (+13.84%)	0.2395 (+4.36%)	0.1478 (+39.96%)	0.0432 (+55.40%)	1.3794 (+37.46%)	8.9044 (+1.81%)
Wikipedia	2.6088 (-24.82%)	0.1892 (+20.05%)	0.2868 (+46.03%)	0.2416 (+169.34%)	0.1046 (+428.28%)	1.6997 (+84.27%)	8.3149 (-2.88%)

(b) Evaluating diversity metrics of synthetic data generation methods from **biomedicine** domain.

Figure 2: Metrics are annotated with ↑ or ↓ arrows which indicate if higher or lower values are better, respectively. **1-GD** refers to 1-Gram diversity and **4-GD** refers to 4-Gram diversity. **MIF** refers to the Mean Inverse Frequency metric (refer to section 4). For a particular domain, diversity metrics for synthetic data generated using template prompting the base LLM are underlined as reference points. We include diversity metrics over a subset of Wikipedia as a generic example of a dataset regarded to be diverse. For each synthetic data generation method and each metric, percentage increases in diversity relative to template prompting are shown in parentheses. Improvements in measured diversity are highlighted in green and reductions in diversity are highlighted in red. All metrics are mean values of 95% CI computed with bootstrap resampling (refer to Appendix D.7). We control for length in all diversity comparisons by constraining synthetic documents to 400 words (Section 3.2) and sampling from a similar-length distribution for other sources (e.g., Common Crawl, Wikipedia; Appendix E).

keyword based which initializes the generation process using random domain-specific keywords synthesized by an agent.

Topic-Aware Seed Selection We introduce a second *topic-aware* seed selection approach using a dynamically adaptive k -NN algorithm. Starting with N seed documents from domain-specific Common Crawl, each assigned an LLM-generated topic label (*Topic Labeling Expert*), we update the seed set every M MetaSynth-generated documents. New seeds are retrieved from the k nearest neighbors of synthesized documents in embedding space, ensuring each has a novel topic label. If insufficient candidates are found, k is incremented³. This ensures topical variation while maintaining semantic relevance to initial seeds. To prevent seed data leakage (see Table 3), MetaSynth always extracts keywords via the *Seed Keyword Extraction Expert*, ensuring synthesis is keyword-driven, regardless of whether

seeds are documents or keywords. Motivated by Eldan and Li (2023), who show that short, diverse, grammatically correct texts (*TinyStories*) induce language learning in small LMs, we cap the length of both synthesized and seed documents at 400 words (approximately 530 tokens). Appendix N contains examples of MetaSynth generated documents.

3.3 MetaSynth-Instruct: Synthetic Instruction Synthesis & Evolution Using Synthetic Documents

We design a meta-prompting driven instruction synthesizer to leverage the synthetic documents synthesized in the previous step (section 3.2) to derive and evolve complex instructions. As part of the meta-prompt, we use a *task description* string to define an instruction as a complex problem about a particular context leveraging various formats and styles e.g. reading comprehension, multiple-choice, fill-in-the-blank and inferential questions. To prevent instruction synthesis from degenerating, the meta-prompt for instruction

³Initially we set $k = 5$; embeddings are computed using <https://huggingface.co/jinaai/jina-embeddings-v2-base-en>

synthesis also contains invocation calls for a certain group of predetermined agents to always be invoked e.g. Document Transformation Expert, Persona Suggestion Expert (inspired by Ge et al. (2024)), Complexity Expert and Question Editor Expert (similar to the suggestor-editor agents proposed by AgentInstruct (Mitra et al., 2024)). In contrast to the work by Xu et al. (2023) and Honovich et al. (2022), in our method the instruction evolution prompts (which involves choosing the method of evolution) are open-ended; composed by the meta-LM taking into account the content of the synthesized document, the responses of agents from previous execution steps and it’s own best judgment. We illustrate an example agentic flow for MetaSynth-*Instruct* in Appendix H and Figure 20. Synthesized instructions (Appendix K, Appendix O) are limited to 100 words and the responses to each instruction are generated using Claude 3 Sonnet with varied prompt formats (see Appendix G).

4 Measuring The Diversity of Generated Synthetic Data

The premise of this work is that diverse data is high quality data. Thus, in lieu of human judgment of diversity; it is necessary to use an appropriate set of automated metrics which can quantify the diversity of LLM generated data such that these measures also align with *human notions of variability and diversity*.

4.1 Metrics:

Task2Vec Diversity Coefficient To quantify semantic and structural diversity in MetaSynth-generated data, we adopt the *Task2Vec* diversity coefficient from Lee et al. (2023). *Task2Vec* formalizes diversity by embedding sampled data batches (e.g., synthesized documents) using the Fisher Information Matrix of a probe network⁴ fine-tuned on the data. The coefficient, defined as the average pairwise cosine distance between *Task2Vec* embeddings (Achille et al., 2019), has been shown to correlate with human diversity judgments.

Compression Ratio & N-Gram Diversity Following the recommendations of Shaib et al. (2024b,a), we select Gzip compression ratio and N-Gram diversity score (ratio of the

⁴We use GPT-2 as the probe network

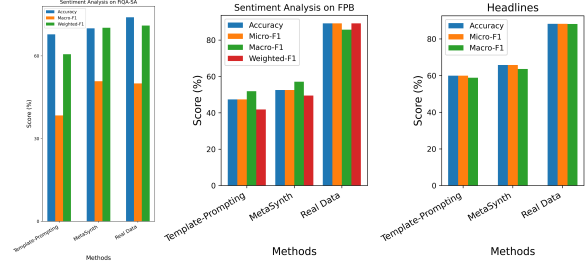


Figure 3: Comparing the performance of BERT fine-tuned on data synthesized with template-prompting and MetaSynth versus real data on: (Left) FiQA-SA; (Middle) FPB; (Right) Headlines.

unique n-gram counts to all n-gram counts in a dataset) as appropriate metrics which can detect aspects of repetition in LLM generated texts (such as the presence of pre-defined syntactic templates).

Remote-Clique & Chamfer Distance

Following (Cox et al. (2021); Li et al. (2023)), we also compute language model embedding based diversity with the *Remote Clique Score* (average mean pairwise distance of a data instance to other instances) and the *Chamfer Distance Score* (average minimum pairwise distance of a data instance to other instances).

Mean Inverse Frequency (MIF) Score We propose an additional metric which captures the average “lexical rarity” of synthesized documents, where instances that use a rarer vocabulary (relative to a reference corpus such as Wikipedia) are assigned high scores, and vice versa, somewhat similar to Inverse Document Frequency from *TF-IDF* (Ramos, 2003). Refer to appendix D for further details on diversity metrics.

Figure 2 shows that MetaSynth documents seeded with Common Crawl are more diverse than those seeded with random keywords, with both exceeding the diversity of template-prompted documents.

5 Experiments and Results

Domain Adaptation Focusing on **continual pre-training** (where the loss is computed on all tokens) and **not supervised instruction fine-tuning** (where the loss is computed only on the response conditioned on the prompt) - we continue to train Mistral-7B (Jiang et al., 2023). As shown in Table 1, 25 million tokens of diverse data synthesized with MetaSynth is sufficient for domain adaptation,

<i>Finance</i>							
CPT Setting	Token Mix	ConvFinQA	NER	FPB	Headline	FiQA_SA	Average
Mistral-7B Base (No CPT)	0M	38.9	58.14	65.09	79.26	75.62	63.40
<u>Real Docs + Template-Prompting Docs</u>	12.5M:12.5M	48.79	52.64	64.24	76.00	74.47	63.23
Real Docs	25M	46.51	55.59	65.07	78.30	76.09	64.31
<u>Real Docs + MetaSynth Docs</u>	12.5M:12.5M	48.59	53.69	67.82	80.14	75.66	65.18
Real Docs + MetaSynth Docs-Instructions-Responses	12.5M:12.5M	43.29	53.77	62.06	79.75	71.57	62.09
Real Docs + MetaSynth Docs-Instructions-Responses	8.33M:16.7M	43.22	52.26	65.66	79.73	72.50	62.67
Real Docs + MetaSynth Instructions-Responses	12.5M:12.5M	47.51	52.08	63.16	79.53	72.56	62.97
Real Docs + MetaSynth Instructions-Responses	8.33M:16.7M	44.43	49.34	63.05	79.68	75.27	62.35
MetaSynth Docs	25M	42.28	48.72	67.37	79.67	73.65	62.34
MetaSynth Docs-Instructions-Responses	25M	49.30	54.64	66.43	83.46	76.13	65.99
<i>Biomedicine</i>							
CPT Setting	Token Mix	PubMedQA	USMLE	MQP	RCT	ChemProt	Average
Mistral-7B (No CPT)	0M	58.20	35.27	67.86	62.55	40.80	52.94
<u>Real Docs + Template-Prompting Docs</u>	12.5M:12.5M	56.40	38.41	67.38	59.80	30.40	50.48
Real Docs	25M	59.70	36.37	62.29	63.70	28.90	50.19
<u>Real Docs + MetaSynth Docs</u>	12.5M:12.5M	60.70	37.31	64.26	67.50	45.00	54.95
Real Docs + MetaSynth Docs-Instructions-Responses	12.5M:12.5M	60.30	37.16	74.75	71.85	38.40	56.49
Real Docs + MetaSynth Docs-Instructions-Responses	8.33M:16.7M	59.50	36.61	76.06	71.05	42.20	57.08
Real Docs + MetaSynth Instructions-Responses	12.5M:12.5M	62.90	35.98	71.80	71.40	39.60	56.34
Real Docs + MetaSynth Instructions-Responses	8.33M:16.7M	60.20	36.44	73.77	71.75	42.10	56.85
MetaSynth Docs	25M	60.20	37.23	70.16	68.15	40.40	55.23
MetaSynth Docs-Instructions-Responses	25M	61.80	36.60	77.87	74.45	50.40	60.22

Table 1: Performance on domain-specific tasks for Mistral-7B under **nine** different continual pre-training (CPT) settings with varying mixing ratios of real and synthetic data. **Bold** indicates the best result for a dataset across all settings within a particular domain. Settings are underlined to indicate the corresponding setting from MetaSynth which can be compared with Template-Prompting.

tested across nine different combinations of mixing Common Crawl texts with synthetic documents and instructions in 1:1 and 1:2 token mixing ratios (refer to appendix C for prompt settings). In Finance, we observe that 25M MetaSynth-generated tokens—without real Common Crawl data—improves the base model by **4.08%** on average, outperforming it on all datasets except NER⁵. A 1:1 mix of real and MetaSynth-generated documents also outperforms the same mix with template-prompted data by **3.08%**. The same holds true for Biomedicine — Continual pretraining on 25M MetaSynth-generated tokens—without real Common Crawl data—also boosts the base model by **13.75%** on average. Similar to finance, a 1:1 real-synthetic document mix outperforms the same mix with template-prompted data by **8.85%**. Overall, in-domain gains are more pronounced in biomedicine, with more types of token mixing ratios improving the base model compared to finance, likely due

⁵This aligns with Cheng et al. (2024a), who note NER’s low benchmark quality, where the base model achieves the highest score

to more specialized terminology and obscure knowledge required for biomedicine, which the base model lacks.

General Evaluation As shown in Table 2, on average, continual pre-training on **MetaSynth generated synthetic data does not compromise the generalizability of the LLM.**

6 Analysis

To evaluate the utility of our instruction synthesizer (MetaSynth-*Instruct*) in creating instructions for more general tasks, we conduct the following analyses:

Creating Data For Fine-tuning Encoders We adapt our instruction synthesizer to generate synthetic data that emulates datasets used in encoder LM evaluation. We modify the task description in the meta-prompt (Appendix K.4) to instruct the meta-LM to generate synthetic training instances resembling each of three finance datasets—Headline News (sarcasm detection), FiQA-SA (aspect-based sentiment analysis), and Financial Phrasebank

Base Model	ARC-ch	ARC-easy	BoolQ	HellaSwag	MMLU	OBQA	PIQA	SIQA	Winogrande	Avg
Mistral-7B	52.1	78.4	82.0	80.4	59.1	44.2	82.3	45.9	73.4	66.4
Finance										
Real Docs + Template Prompting Docs	53.8	78.4	78.0	80.7	59.0	45.6	81.9	48.1	71.4	66.3
Real Docs + MetaSynth Docs	55.9	77.3	84.3	80.7	58.5	44.4	81.1	49.6	71.7	67.1
MetaSynth Docs-Instr-Responses	50.9	75.1	84.1	79.4	56.3	43.0	80.7	48.1	69.3	65.2
Biomedicine										
Real Docs + Template Prompting Docs	54.9	79.5	80.8	81.1	58.1	45.4	82.6	46.9	71.7	66.8
Real Docs + MetaSynth Docs	53.4	76.1	83.5	80.6	58.0	44.6	81.0	46.9	70.2	66.0
MetaSynth Docs-Instr-Responses	54.2	75.2	83.2	79.1	57.5	43.2	81.0	47.5	70.8	65.8

Table 2: General evaluation across domains and Settings. Real docs + Template Prompting docs refers to Continual Pre-training (CPT) over 12.5M tokens of synthetic data generated with template prompting mixed with 12.5 tokens of Common Crawl data. Real Docs + MetaSynth Docs refers to CPT over 12.5M tokens of synthetic data generated by our method mixed with 12.5M tokens of Common Crawl data. MetaSynth Docs-Instr-Responses refers to CPT over 25M tokens of MetaSynth documents and their associated synthetic instruction-response pairs.

Dataset	EM-1	EM-2	EM-3	EM-5	EM-10
General Datasets					
ConvFinQA	0.9784	0.7756	0.2603	0.0310	0.0000
NER	0.9923	0.7416	0.2431	0.0287	0.0000
FPB	0.9681	0.7024	0.3137	0.0222	0.0000
Headline	0.9957	0.6752	0.1727	0.0075	0.0000
FiQA_SA	0.9619	0.5745	0.1852	0.0069	0.0000
Biomedical Datasets					
ChemProt	0.9329	0.5933	0.2298	0.0111	0.0000
MQP	0.9893	0.8411	0.4211	0.0279	0.0000
PubMedQA	0.9867	0.7431	0.3214	0.0257	0.0000
RCT	0.9886	0.7726	0.4108	0.0422	0.0000
USMLE	0.9951	0.8137	0.4190	0.0495	0.0000

Table 3: Data contamination check results. EM-N stands for Exact Match N -gram overlap as a substring between the reference texts from each benchmark dataset and potentially contaminated target texts from Common Crawl (Real Docs).

(sentiment analysis)—selected for their simplicity and prior use in Li et al. (2023)’s work. For each dataset, we generate a small set of synthetic instances with both MetaSynth and template-prompting using 3-shot examples. Fine-tuning a BERT-based classifier on the generated data and evaluating it on the real test partition of each dataset shows that models trained on MetaSynth data outperform those trained on data synthesized by template prompting but remain behind those fine-tuned on real data, consistent with Li et al. (2023)’s findings (Figure 3).

Instruction-Response Quality Following Cheng et al. (2024a), we also analyze our synthesized instruction-response pairs in terms of *context relevance*, *response accuracy*, *task diversity* and *win rate*. Evaluating 1000 sampled instruction-response pairs from each domain and using Claude 3 Opus (Anthropic,

2024) as a judge. Table 4, Appendix F shows that our synthesized instruction-response pairs for finance exhibit greater task diversity and slightly higher relevance and accuracy scores than Biomedicine, yet 25M biomedical tokens still yield greater improvements to the base model, suggesting that achieving comparable gains in finance would require substantially more data than what we synthesized due to it being a more generic domain. For both domains, a Mistral-7B model continually pre-trained on 25M MetaSynth tokens also attains higher win rates against Claude 3 Sonnet relative to the base model (Figure 11, appendix F). Appendix D.8 shows that instruction diversity (computed using our metrics) relative to *Instruction-Pretraining* (Cheng et al., 2024a) is lower for MetaSynth, attributable to using 1B tokens of real text to synthesize instructions versus our 25M synthetic tokens.

Is it Data Contamination? We assess cross-contamination between Common Crawl (*Real Docs*) and domain-specific benchmarks e.g., ConvFinQA using a 10-gram substring match method (Ben Allal et al., 2024; OpenAI et al., 2024), deeming an example contaminated if a substring appears in *Real Docs*. Table 3 shows no contamination between our selected Common Crawl seeds and evaluation datasets.

7 Conclusion

We propose METASYNTH, a method that leverages meta-prompting and agentic scaffolding to generate diverse documents and instructions. We demonstrate its efficacy by synthesizing diverse data and then continually pre-training Mistral-7B on it, yielding significant improvements in two domains, without degrading the model on general tasks.

Limitations

Our work has several limitations worth noting. First, our approach of iteratively refining each synthetic instance to be more diverse, while keeping track of all previously generated instances incurs a significant inference cost when synthesizing a collection of documents. A run-time analysis of our method reveals it takes approximately approximately 3.6 minutes to synthesize a single document (or 3 hours to generate 50 documents starting from initial seeds). While this inference-time trade off is intentional (with the objective of increasing the diversity of generated data), it is still an important consideration, especially when operating in resource-constrained settings; for e.g., template-prompting based methods can synthesize a document or instruction in just a few seconds.

A significant challenge also lies in the stability of our agentic workflow. We observe that our procedure is prone to breakdowns, requiring many iterations to be discarded. This instability suggests that more robust methods for maintaining coherent meta-level control may be needed for deploying our approach practically.

Despite our efforts to generate diverse synthetic data, an analysis of the diversity of our synthesized instructions reveals that it still falls short compared to the diversity of instructions evolved from real corpora (as shown in Appendix D.8), thus synthesized documents are still not a true replacement for human generated documents (which is expected).

Furthermore, we find that automatic evaluation metrics for assessing data diversity may not always align well with human judgments. A concrete example of this emerges in our finance domain experiments, where we observe strong biases in the generated content towards specific topics like “ESG”, “DeFi”, and “cryptocurrency”. These biases likely stem from the underlying LLM–Claude 3 Sonnet’s–post-training alignment. This highlights a broader challenge with synthetic data generation methods: ensuring that the generated data not only appears diverse by automated metrics but also maintains domain-appropriate distributions and high diversity by human standards. We leave a human evaluation of the diversity of data synthesized by MetaSynth as future work.

References

- Alessandro Achille, Michael Lam, Rahul Tewari, Avinash Ravichandran, Subhansu Maji, Charles Fowlkes, Stefano Soatto, and Pietro Perona. 2019. [Task2vec: Task embedding for meta-learning](#). *Preprint*, arXiv:1902.03545.
- Anthropic. 2024. [Introducing the next generation of claude](#). Accessed: 2025-01-25.
- Loubna Ben Allal, Anton Lozhkov, Guilherme Penedo, Thomas Wolf, and Leandro von Werra. 2024. [Cosmopedia](#).
- Jiuhai Chen, Rifaa Qadri, Yuxin Wen, Neel Jain, John Kirchenbauer, Tianyi Zhou, and Tom Goldstein. 2024. [Genqa: Generating millions of instructions from a handful of prompts](#). *Preprint*, arXiv:2406.10323.
- Zhiyu Chen, Shiyang Li, Charese Smiley, Zhiqiang Ma, Sameena Shah, and William Yang Wang. 2022. [Convfinqa: Exploring the chain of numerical reasoning in conversational finance question answering](#). *Preprint*, arXiv:2210.03849.
- Daixuan Cheng, Yuxian Gu, Shaohan Huang, Junyu Bi, Minlie Huang, and Furu Wei. 2024a. [Instruction pre-training: Language models are supervised multitask learners](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2529–2550, Miami, Florida, USA. Association for Computational Linguistics.
- Daixuan Cheng, Shaohan Huang, and Furu Wei. 2024b. [Adapting large language models to domains via reading comprehension](#). *Preprint*, arXiv:2309.09530.
- Samuel Rhys Cox, Yunlong Wang, Ashraf Abdul, Christian von der Weth, and Brian Y. Lim. 2021. [Directed diversity: Leveraging language embedding distances for collective creativity in crowd ideation](#). In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, CHI ’21, New York, NY, USA. Association for Computing Machinery.
- Franck Dernoncourt and Ji Young Lee. 2017. [PubMed 200k RCT: a dataset for sequential sentence classification in medical abstracts](#). In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 308–313, Taipei, Taiwan. Asian Federation of Natural Language Processing.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#). *Preprint*, arXiv:1810.04805.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong

- Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*.
- Ronen Eldan and Yuanzhi Li. 2023. [Tinystories: How small can language models be and still speak coherent english?](#) *Preprint*, arXiv:2305.07759.
- 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.
- Matthias Gerstgrasser, Rylan Schaeffer, Apratim Dey, Rafael Rafailov, Henry Sleight, John Hughes, Tomasz Korbak, Rajashree Agrawal, Dhruv Pai, Andrey Gromov, Daniel A. Roberts, Diyi Yang, David L. Donoho, and Sanmi Koyejo. 2024. [Is model collapse inevitable? breaking the curse of recursion by accumulating real and synthetic data](#). *Preprint*, arXiv:2404.01413.
- Aaron Grattafiori et al. 2024. [The llama 3 herd of models](#).
- Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. [Large language model based multi-agents: A survey of progress and challenges](#). In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24*, pages 8048–8057. International Joint Conferences on Artificial Intelligence Organization. Survey Track.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. [Training compute-optimal large language models](#). *Preprint*, arXiv:2203.15556.
- Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. 2022. [Unnatural instructions: Tuning language models with \(almost\) no human labor](#). *Preprint*, arXiv:2212.09689.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *Preprint*, arXiv:2310.06825.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. [What disease does this patient have? a large-scale open domain question answering dataset from medical exams](#). *Preprint*, arXiv:2009.13081.
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. 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 2567–2577.
- Seungone Kim, Juyoung Suk, Xiang Yue, Vijay Viswanathan, Seongyun Lee, Yizhong Wang, Kiril Gashteovski, Carolin Lawrence, Sean Welleck, and Graham Neubig. 2024. [Evaluating language models as synthetic data generators](#). *Preprint*, arXiv:2412.03679.
- Jens Vindahl Kringelum, Sonny Kim Kj  rulf, S  ren Brunak, Ole Lund, Tudor I. Oprea, and Olivier Taboureau. 2016. [Chemprot-3.0: a global chemical biology diseases mapping](#). *Database: The Journal of Biological Databases and Curation*, 2016.
- Alycia Lee, Brando Miranda, and Sanmi Koyejo. 2023. Beyond scale: the diversity coefficient as a data quality metric demonstrates llms are pre-trained on formally diverse data. *arXiv preprint arXiv:2306.13840*.
- Junlong Li, Jinyuan Wang, Zhuosheng Zhang, and Hai Zhao. 2024. [Self-prompting large language models for zero-shot open-domain QA](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 296–310, Mexico City, Mexico. Association for Computational Linguistics.
- Zhuoyan Li, Hangxiao Zhu, Zhuoran Lu, and Ming Yin. 2023. [Synthetic data generation with large language models for text classification: Potential and limitations](#). *Preprint*, arXiv:2310.07849.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pre-training approach](#). *Preprint*, arXiv:1907.11692.
- Macedo Maia, Siegfried Handschuh, Andr   Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. [Www’18 open challenge: Financial opinion mining and question answering](#). In *Companion Proceedings of the The Web Conference 2018, WWW ’18*, page 1941–1942, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.
- Pekka Malo, Ankur Sinha, Pyry Takala, Pekka Korhonen, and Jyrki Wallenius. 2013. [Good debt or bad debt: Detecting semantic orientations in economic texts](#). *Preprint*, arXiv:1307.5336.

800	Clara H. McCreery, Namit Katariya, Anitha Kannan, Manish Chablani, and Xavier Amatriain. 2020. Effective transfer learning for identifying similar questions: Matching user questions to covid-19 faqs . In <i>Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining</i> , KDD '20, page 3458–3465, New York, NY, USA. Association for Computing Machinery.	
801		
802		
803		
804		
805		
806		
807		
808		
809	Arindam Mitra, Luciano Del Corro, Guoqing Zheng, Shweti Mahajan, Dany Rouhana, Andres Codas, Yadong Lu, Wei ge Chen, Olga Vrousos, Corby Rosset, Fillipe Silva, Hamed Khanpour, Yash Lara, and Ahmed Awadallah. 2024. Agentinstruct: Toward generative teaching with agentic flows . <i>Preprint</i> , arXiv:2407.03502.	
810		
811		
812		
813		
814		
815		
816	Sania Nayab, Giulio Rossolini, Marco Simoni, Andrea Saracino, Giorgio Buttazzo, Nicolamaria Manes, and Fabrizio Giacomelli. 2025. Concise thoughts: Impact of output length on llm reasoning and cost . <i>Preprint</i> , arXiv:2407.19825.	
817		
818		
819		
820		
821	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2024. Gpt-4 technical report . <i>arXiv preprint arXiv:2303.08774</i> .	
822		
823		
824		
825		
826		
827	Juan Enrique Ramos. 2003. Using tf-idf to determine word relevance in document queries .	
828		
829	Julio Cesar Salinas Alvarado, Karin Verspoor, and Timothy Baldwin. 2015. Domain adaption of named entity recognition to support credit risk assessment . In <i>Proceedings of the Australasian Language Technology Association Workshop 2015</i> , pages 84–90, Parramatta, Australia.	
830		
831		
832		
833		
834		
835	Mohamed El Amine Seddik, Swei-Wen Chen, Soufiane Hayou, Pierre Youssef, and Merouane Debah. 2024. How bad is training on synthetic data? a statistical analysis of language model collapse . <i>Preprint</i> , arXiv:2404.05090.	
836		
837		
838		
839		
840	Chantal Shaib, Joe Barrow, Jiuding Sun, Alexa F. Siu, Byron C. Wallace, and Ani Nenkova. 2024a. Standardizing the measurement of text diversity: A tool and a comparative analysis of scores . <i>Preprint</i> , arXiv:2403.00553.	
841		
842		
843		
844		
845	Chantal Shaib, Yanai Elazar, Junyi Jessy Li, and Byron C. Wallace. 2024b. Detection and measurement of syntactic templates in generated text . <i>Preprint</i> , arXiv:2407.00211.	
846		
847		
848		
849	Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. 2024. The curse of recursion: Training on generated data makes models forget . <i>Preprint</i> , arXiv:2305.17493.	
850		
851		
852		
853		
854	Ankur Sinha and Tanmay Khandait. 2020. Impact of news on the commodity market: Dataset and results . <i>Preprint</i> , arXiv:2009.04202.	
855		
856		
	Mirac Suzgun and Adam Tauman Kalai. 2024. Meta-prompting: Enhancing language models with task-agnostic scaffolding .	857
		858
		859
	Pablo Villalobos, Anson Ho, Jaime Sevilla, Tamay Besiroglu, Lennart Heim, and Marius Hobbhahn. 2024. Will we run out of data? limits of llm scaling based on human-generated data . <i>Preprint</i> , arXiv:2211.04325.	860
		861
		862
		863
		864
	Fanqi Wan, Xinting Huang, Tao Yang, Xiaojun Quan, Wei Bi, and Shuming Shi. 2023. Explore-instruct: Enhancing domain-specific instruction coverage through active exploration . <i>Preprint</i> , arXiv:2310.09168.	865
		866
		867
		868
		869
	Yizhong Wang, Swaroop Mishra, Pegah Alipoor-molabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Maitreya Patel, Kuntal Kumar Pal, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Shailaja Keyur Sampat, Savan Doshi, Siddhartha Mishra, Sujan Reddy, Sumanta Patro, Tanay Dixit, Xudong Shen, Chitta Baral, Yejin Choi, Noah A. Smith, Hannaneh Hajishirzi, and Daniel Khoshnab. 2022. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks . <i>Preprint</i> , arXiv:2204.07705.	870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models . <i>Preprint</i> , arXiv:2201.11903.	888
		889
		890
		891
		892
	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryan W White, Doug Burger, and Chi Wang. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation . <i>Preprint</i> , arXiv:2308.08155.	893
		894
		895
		896
		897
		898
		899
	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhao Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions . <i>Preprint</i> , arXiv:2304.12244.	900
		901
		902
		903
		904
	Ran Xu, Hejie Cui, Yue Yu, Xuan Kan, Wenqi Shi, Yuchen Zhuang, Wei Jin, Joyce Ho, and Carl Yang. 2025. Knowledge-infused prompting: Assessing and advancing clinical text data generation with large language models . <i>Preprint</i> , arXiv:2311.00287.	905
		906
		907
		908
		909
		910
	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models . <i>Preprint</i> , arXiv:2210.03629.	911
		912
		913
		914

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.

Yifan Zhang, Yang Yuan, and Andrew Chi-Chih Yao. 2024. [Meta prompting for ai systems](#). *Preprint*, arXiv:2311.11482.

A Related Work

Prior work on generating synthetic data with LLMs has primarily focused on post-training data synthesis, particularly for conversational data or instructions (Honovich et al., 2022; Xu et al., 2023; Chen et al., 2024; Ding et al., 2023) inter alia. Recent work like AgentInstruct uses predefined taxonomies and agentic workflows to generate instruction-response pairs from real corpora. While similar to our approach in using iterative refinement, our method differs by leveraging the meta-model’s reasoning to dynamically select synthesis flows rather than sampling from fixed taxonomies. Unlike post-training approaches that compute loss only on responses, our method aligns with instruction pre-training approaches (Cheng et al., 2024b,a) by computing loss on both prompts and responses. However, we operate with significantly smaller token counts - using 26.5M tokens from domain-specific Common Crawl splits (approximately 50K documents) compared to AdaptLLM’s (Cheng et al., 2024b) use of billions of tokens from real corpora (5.4B medical, 1.2B finance). Additionally, our method is unsupervised, generating instructions from synthetic texts, whereas Instruction-Pretraining (Cheng et al., 2024a) leverages an instruction synthesizer trained on at least 1B tokens of real corpora. Another approach PersonaHub (Ge et al., 2024) first samples 1 Billion persona’s from 10^{14} tokens of web scale text and then uses a template prompt e.g., “create data with persona” to synthesize instances. Given the large scale of real data used in PersonaHub it is not comparable to our method.

B Meta-Prompting

B.1 Algorithmic Procedure

Let \mathbb{S} be the set of finite strings, with \emptyset denoting the empty string. A test-time query $x \in \mathbb{S}$ represents a natural language task. The fixed

Algorithm 1 Meta Prompting

Require: $\text{LM} : \mathcal{S} \rightarrow \mathcal{S}$, $x, \text{error} \in \mathcal{S}$; $T \in \mathbb{N}$;
 $t_{\text{init}}, t_{\text{mid}}, t_{\text{exp}}, e_{\text{exp}}, e_{\text{ret}} : \mathcal{S} \rightarrow \mathcal{S}$

- 1: $\mathcal{H}_1 \leftarrow t_{\text{init}}(x)$
- 2: **for** $t \in [1, \dots, T]$ **do**
- 3: $y_t \leftarrow \text{LM}(\mathcal{H}_t)$
- 4: **if** $e_{\text{exp}}(y_t) \neq \emptyset$ **then**
- 5: $\text{prompt} \leftarrow t_{\text{exp}}(e_{\text{exp}}(y_t))$
- 6: $z_t \leftarrow \text{LM}(\text{prompt})$
- 7: $\mathcal{H}_{t+1} \leftarrow \mathcal{H}_t \oplus t_{\text{mid}}(z_t)$ {Meta Model provided expert instructions}
- 8: **else if** $e_{\text{ret}}(y_t) \neq \emptyset$ **then**
- 9: **return** $e_{\text{ret}}(y_t)$ {Meta Model returned end of generation token}
- 10: **else**
- 11: $\mathcal{H}_{t+1} \leftarrow \mathcal{H}_t \oplus \text{error}$ {Meta Model formatting error}
- 12: **end if**
- 13: **end for**

language model LM operates from \mathbb{S} to \mathbb{S} , taking a prompt history \mathcal{H} as input and producing an output. Template functions t_{init} , t_{mid} , and t_{exp} map \mathbb{S} to \mathbb{S} , formatting input/output for the meta-LM and agent/expert models. String extractors e_{exp} and e_{ret} retrieve substrings enclosed by delimiters, while \oplus denotes string concatenation, and $\text{error} \in \mathbb{S}$ represents error messages. At each iteration, \mathcal{H}_t guides LM to either return a response or consult an agent, with instructions extracted via e_{exp} . Agents process only what is shared with them by the meta-LM, and their outputs are formatted with t_{mid} . A final response is extracted using e_{ret} and returned. If neither a final response nor a call to an agent is made, error is appended to \mathcal{H}_t for error handling.

B.2 Conditional Instance Generation

The idea of conditional instance generation as applied to synthesizing documents using seed keywords is expressed in Algorithm 2. After each synthesized document, the set of seed keywords is expanded with related yet distinct terms. Each synthesized document must satisfy the following two criteria: conform to the current seed set, and be distinct from all previously synthesized documents. This is achieved by using a meta-LLM to continuously categorize and keep track of summaries of all prior documents.

Algorithm 2 Conditional Instance Generation

Require: S_0 : initial set of seeds;
 θ : parameters of data generating LM;
 $\text{div}(\cdot, \cdot)$: implicit diversity measure between a set of instances; $T \in \mathbb{N}$

- 1: **Step 1** Generate an initial instance: $I_0 \sim p(\cdot | S_0; \theta)$
- 2: **Step 2** Expand seed set and then generate another instance:
- 3: $S_1 = \text{ExpandSeeds}(S_0, \{I_0, I_1\})$
- 4: $I_1 = \arg \max_I \left[p(I | I_0, S_1; \theta) \times \mathbb{E}[\text{div}(I_0, I)] \right]$
 \therefore subject to I conforming to S_1
- 5: **Step 3** Iteratively generate additional instances:
- 6: **for** $i = 2$ to T **do**
- 7: $S_i = \text{ExpandSeeds}(S_{i-1}, \{I_0, \dots, I_i\})$
 $I_i = \arg \max_I \left[p(I | \{I_0, \dots, I_{i-1}\}, S_i; \theta) \right.$
 $\quad \times \mathbb{E}[\text{div}(\{I_0, \dots, I_{i-1}\}, I)] \left. \right]$
 \therefore subject to I conforming to S_i
- 8: **end for**
- 9: **return** $\{I_0, \dots, I_T\}$

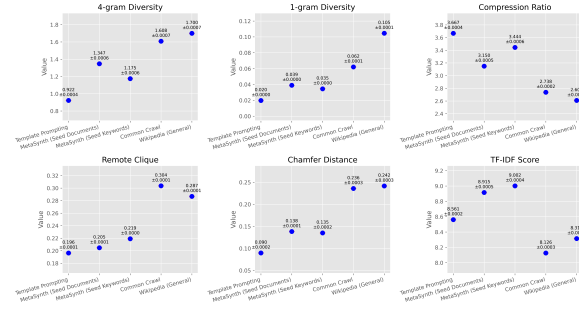


Figure 4: Distribution of diversity metrics for documents synthesized by MetaSynth versus other types of documents (e.g., those generated with template-prompting or real data).

C Prompt Settings & Datasets For Domain Adaptation Experiments

We follow the prompting settings of AdaptLLM (Cheng et al., 2024b): for biomedicine domain, we evaluate zero-shot performance on PubMedQA (Jin et al., 2019) and USMLE (Jin et al., 2020), few-shot performance on ChemProt (Kringelum et al., 2016), MQP (McCreery et al., 2020) and RCT (Dernoncourt and Lee, 2017); for finance domain, we evaluate zero-shot performance on ConvFinQA (Chen et al., 2022) and few-shot performance on FPB (Malo et al., 2013), FiQA SA (Maia et al., 2018), Headline (Sinha and Khandait, 2020), and NER (Salinas Alvarado et al., 2015).

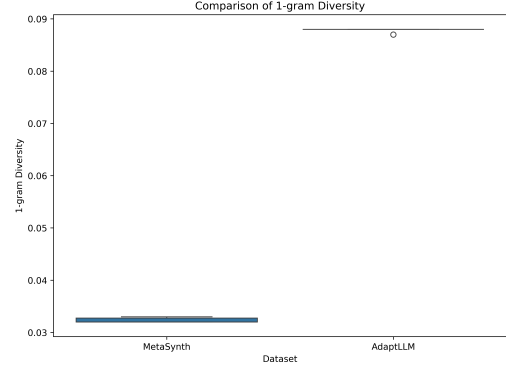


Figure 5: Distribution of 1-Gram diversity between instructions synthesized by MetaSynth-Instruct versus Instruction-Pretraining.

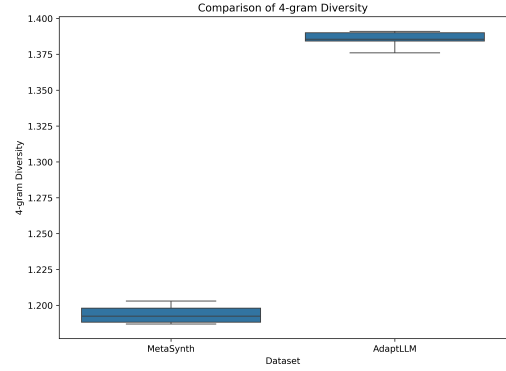


Figure 6: Distribution of 4-Gram diversity between instructions synthesized by MetaSynth-Instruct versus Instruction-Pretraining.

D Diversity Metrics

D.1 Task2Vec Diversity Coefficient

The Task2Vec diversity coefficient proposed by Lee et al. (2023) quantifies the intrinsic variability of a dataset by measuring the distinctness of different data batches, which can be measured for each batch by computing the diagonal of the Fisher Information Matrix (FIM) using a fixed GPT-2 probe network. Intuitively, if a dataset is rich in latent concepts, different batches will fine-tune the final layer of GPT-2 in diverse ways, resulting in Task2Vec embeddings that are more dissimilar (i.e., have larger pairwise cosine distances). Thus, a dataset containing a wide variety of topics and styles should exhibit a higher diversity coefficient than a more homogeneous dataset. The coefficient is calculated as follows:

- **Sampling Batches:**

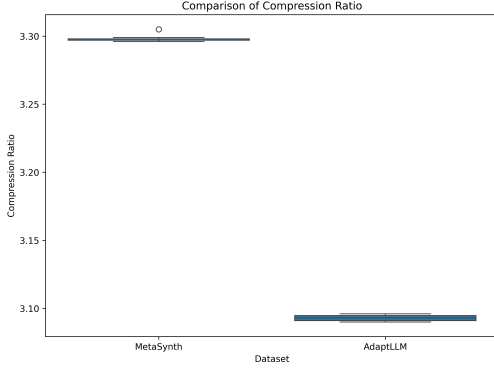


Figure 7: **Distribution of Compression Ratios** between instructions synthesized by MetaSynth-*Instruct* versus Instruction-Pretraining.

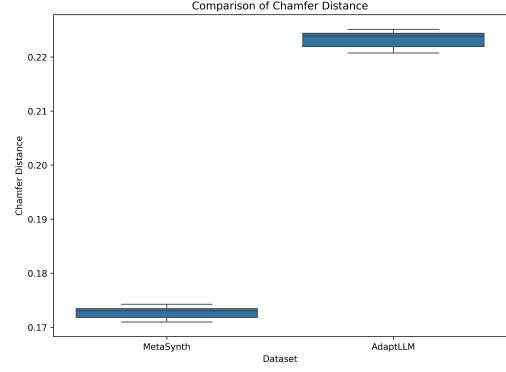


Figure 9: **Distribution of Chamfer Distance** between instructions synthesized by MetaSynth-*Instruct* versus Instruction-Pretraining.

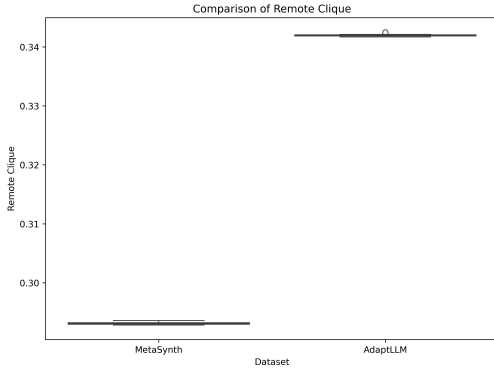


Figure 8: **Distribution of Remote Clique Distance** between instructions synthesized by MetaSynth-*Instruct* versus Instruction-Pretraining.

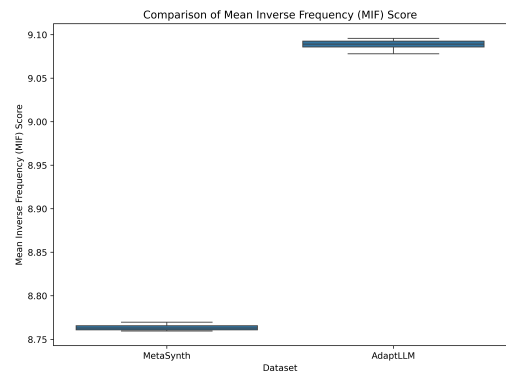


Figure 10: **Distribution of Mean Inverse Frequency (MIF) Score** between instructions synthesized by MetaSynth-*Instruct* versus Instruction-Pretraining.

Sample M batches from a dataset D e.g., the corpus of documents synthesized with MetaSynth. Each batch B_i (for $i = 1, \dots, M$) consists of n text sequences:

$$B_i = \{x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}\}.$$

- **Fine-Tuning the Probe Network:**

For each batch B_i , we fine-tune the final layer of the fixed GPT-2 probe network f_w using a next-token prediction objective. All layers except the final one remain frozen.

- **Computing Gradients:**

For each sequence $x \in B_i$ and each token position t , we then compute the gradient of the log-likelihood with respect to the final-layer parameters:

$$g_t^{(i)} = \nabla_w \log \hat{p}_w(x_t | x_{1:t-1}).$$

- **Estimating the Fisher Information Matrix (FIM):**

For each batch B_i , FIM is approximated by taking the expected outer product of the gradients:

$$\hat{F}_{B_i} = \mathbb{E}_{(x,t) \sim B_i} [g_t^{(i)} (g_t^{(i)})^\top].$$

- **Extracting the Task2Vec Embedding:**

The Task2Vec embedding f_{B_i} for each batch B_i is defined as the diagonal of the FIM:

$$f_{B_i} = \text{diag}(\hat{F}_{B_i}).$$

- **Computing Pairwise Cosine Distances:**

For every distinct pair of batches (B_i, B_j) with $i < j$, we then compute the cosine

distance between their embeddings:

$$d_{ij} = d(f_{B_i}, f_{B_j}).$$

- **Calculating the Diversity Coefficient:** The Task2Vec diversity coefficient is estimated as the average pairwise cosine distance across all batches:

$$\hat{div}(D) = \frac{2}{M(M-1)} \sum_{1 \leq i < j \leq M} d_{ij}.$$

D.2 Compression Ratio

Compression Ratio (CR) Text compression algorithms identify redundancy in variable-length sequences:

$$CR(D) = \frac{\text{size of } D \oplus}{\text{compressed size of } D \oplus} \quad (1)$$

where $D \oplus$ denotes the dataset D concatenated into a single string.

D.3 N-Gram Diversity

N-Gram Diversity Score (NGD) NGD extends the idea of token-type ratio (i.e., the unique token count divided by the total count of tokens) to longer n -grams:

$$NGD(D) = \sum_{n=1}^4 \frac{\# \text{ unique } n\text{-grams in } D \oplus}{\# n\text{-grams in } D \oplus} \quad (2)$$

where $D \oplus$ denotes the dataset D concatenated into a single string.

D.4 Remote Clique

Remote-Clique Distance Average of mean pairwise distances:

$$\frac{1}{N^2} \sum_{i,j} d(\mathbf{x}_i, \mathbf{x}_j) \quad (3)$$

where \mathbf{x}_i represents a document embedding vector computed by a language model.

D.5 Chamfer Distance

Chamfer Distance Average of minimum pairwise distances, also computed over document embeddings:

$$\frac{1}{N} \sum_{i=1}^N \min_{j \neq i} d(\mathbf{x}_i, \mathbf{x}_j) \quad (4)$$

D.6 Mean Inverse Frequency (MIF)

This metric captures the use of rare vocabulary in synthesized documents. For each word, we calculate its inverse frequency value based on a reference corpus (in this case Wikipedia). We then average these values over all words in the document to produce a document-level score that captures lexical rarity.

D.7 Diversity Distribution for MetaSynth Documents Vs Real Documents

Figure 4 presents diversity metrics computed over 5000 synthetic documents with 95% confidence intervals via 1000 bootstrap resamples. MetaSynth documents generated using seed documents exhibit consistently higher diversity and greater similarity to real data, as measured by these metrics, compared to those generated from seed keywords. In turn, MetaSynth even with seed keywords produces more diverse outputs than template-prompting (which uses seed documents as in-context exemplars).

D.8 Diversity Distribution for MetaSynth Instructions Vs Instruction-Pretraining

Figures 5, 6, 7, 8, 9, 10 illustrate the variance for diversity metrics between MetaSynth-*Instruct* and *Instruction-Pretraining* (Cheng et al., 2024a). MetaSynth instructions exhibit lower diversity, as they evolve solely from synthetic documents, whereas Instruction-Pretraining leverages a 1B-token real corpus.

	Accuracy	Relevance	# Category
BioMed.	82.0	91.0	16
Finance	83.0	93.0	23

Table 4: **Response accuracy, context relevance, and number of task categories** of the instruction-response pairs synthesized by MetaSynth.

E Length Distribution: MetaSynth Documents Vs Real Documents

Figures 14, 15, 16, 17, 18, 19 show the length distributions of each type of synthetic or real document used in our work.

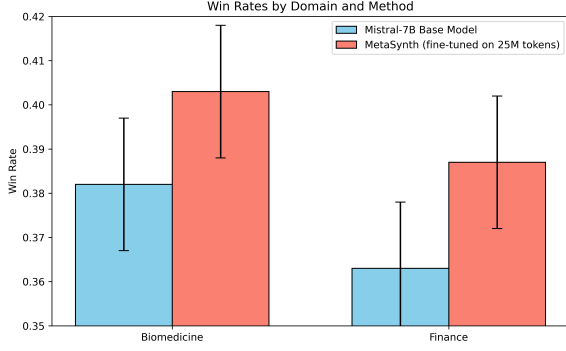


Figure 11: **Against Claude 3 Sonnet (data generating LLM)**: Win-rates shown for Mistral-7B pretrained on 25M tokens of MetaSynth Documents-Instructions-Responses versus the non-pretrained base model (judged by Claude 3 Opus).

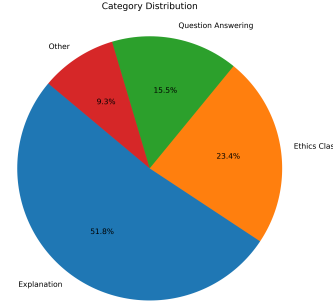


Figure 13: **Distribution of task scenarios synthesized by MetaSynth-Instruct** in instruction-response pairs from Biomedicine domain.

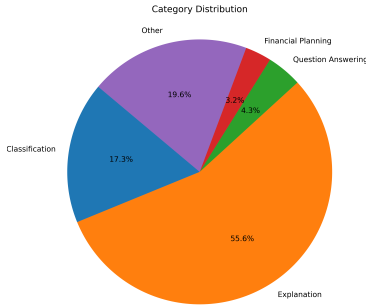


Figure 12: **Distribution of task scenarios synthesized by MetaSynth-Instruct** in instruction-response pairs from Finance domain.

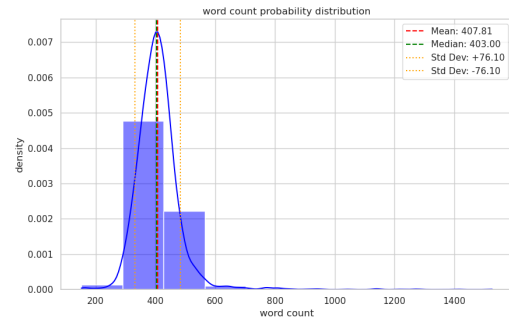


Figure 14: Length distribution (in word count) of documents synthesized by *MetaSynth* from the Finance domain.

F MetaSynth-Instruct Instruction-Response Analysis

Figure 12 and 13 show the percentages of task scenarios from Wang et al. (2022) that occur in a sample of instruction-response pairs synthesized by MetaSynth, for each domain. Table 4 shows the number of unique task scenarios that occur in this sample, along with response accuracy and context relevance.

Response Accuracy Claude 3 Opus (Anthropic, 2024) is prompted to assess whether a response is accurate based on the instruction and context. A binary indicator score is used to compute accuracy.

Context Relevance The same LLM is also prompted to judge whether the instruction synthesized by MetaSynth is relevant to the context (a synthetic document) given the synthesized response to the context-instruction pair.

Win Rate Win-rates against the syn-

thetic data-generating LLM (Claude 3 Sonnet)—which also synthesized responses to its own instructions—are evaluated using Claude 3 Opus as the judge (figure 11). Models continually pre-trained on MetaSynth-generated synthetic data achieve higher win rates than their respective base models. Specifically, Mistral-7B continually pretrained on 25M MetaSynth tokens achieves a 40.3% win rate against Claude 3 Sonnet in biomedicine, outperforming the base (non-pretrained) model’s 38.2%. In finance, it wins 38.7% of the time compared to the base model’s 36.3%, indicating the utility of our synthetic data.

G Templates for Synthesizing Responses to MetaSynth Instructions

To further elicit diverse responses to our synthesized instructions, we reformat each context and its associated instruction pairs through templated variations. Specifically, for each

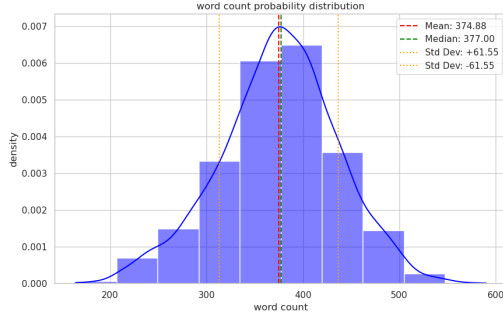


Figure 15: Length distribution (in word count) of documents synthesized by Template-Prompting from the Finance domain.

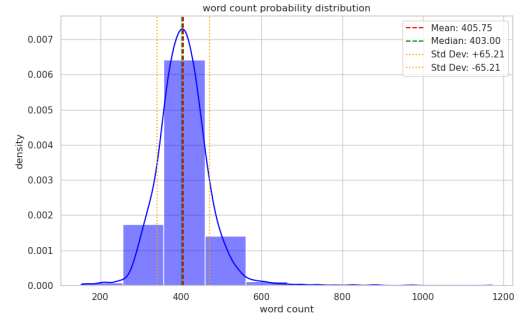


Figure 17: Length distribution (in word count) of documents synthesized by *MetaSynth* from the Biomedicine domain.

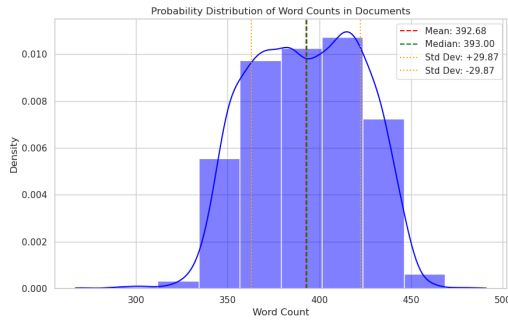


Figure 16: Length distribution (in word count) of Common Crawl documents from the Finance domain.

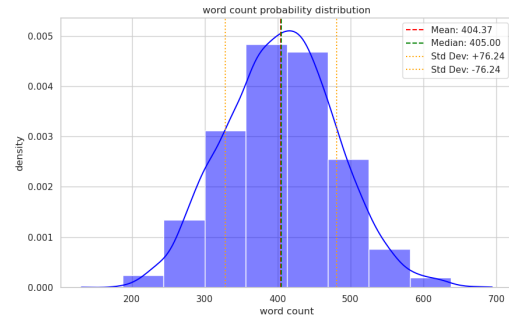


Figure 18: Length distribution (in word count) of documents synthesized by Template-Prompting from the Biomedicine domain.

pair, we apply one of three formats: free-form completion, chain-of-thought (CoT) completion (Wei et al., 2023), and constrained chain-of-thought (cCoT) completion (Nayab et al., 2025). In the cCoT case, a random word limit (a multiple of 50 between 50 and 500) is inserted into the template. We then construct multiple prompt variants by concatenating each context with randomly sampled subsets of these reformatted instructions—ensuring that every example is incorporated at least once—until a full set of variations is obtained. From these variations, we subsequently sample a diverse, non-redundant set of context–instruction pairs for response synthesis.

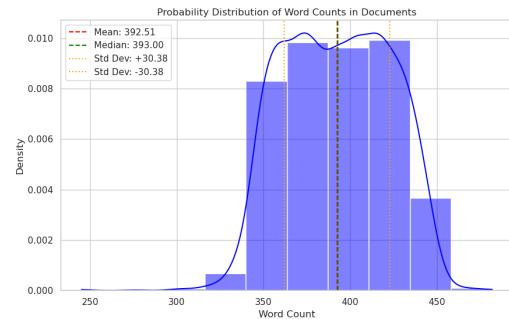


Figure 19: Length distribution (in word count) of Common Crawl documents from the Biomedicine domain.

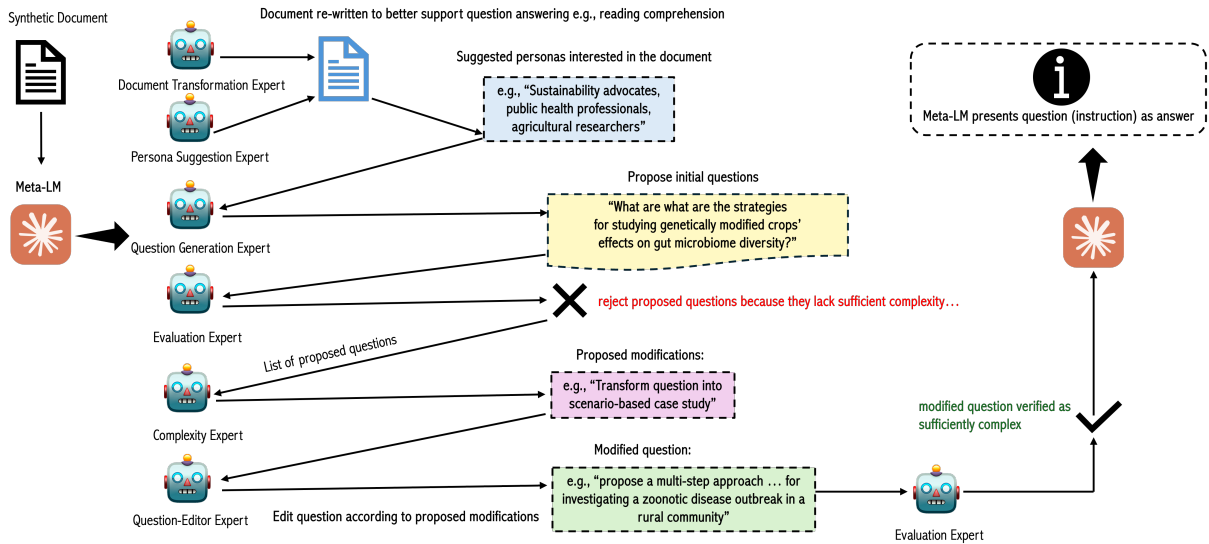


Figure 20: Demonstration of an example METASYNTH-Instruct agentic workflow for synthesizing an instruction from a synthetic biomedical document. A meta-LM orchestrates various expert agents that iteratively refine and generate complex instructions in the form of questions conditioned on the text of the synthetic document.

H MetaSynth-Instruct

Using the above figure as reference, we describe a possible execution history for synthesizing an instruction from a biomedical document as follows:

(1) Given a synthetic document on agriculture and its impact on human health, the meta-LM conjures the following experts to consult: Document Transformation Expert, Persona Suggestion Expert, Question Generation, Evaluation Expert, Complexity Expert, Question Editor expert and two Domain Experts in healthcare and agriculture.

(2) Given the original document text, the Document Transformation Expert is first called by the meta-LM, who identifies the mention of Genetically Modified Organisms and risks/controversies surrounding their use. This expert then reformulates the document to focus more on this aspect.

(3) The meta-LM then calls on the two domain experts (healthcare & agriculture) to analyze both the original and rewritten document(s) and to provide foundational knowledge. At the same time, the meta-LM also calls Persona Suggestion Expert for a list of readers who would find the document engaging; this expert suggests that sustainability advocates, public health professionals and agricultural researchers would be interested in reading the document.

(4) This feedback along with other information aggregated from various experts by the meta-LM is then passed to a Question Generation Expert which then proposes a set of initial questions e.g. “what are the strategies for studying genetically modified crops’ effects on gut microbiome diversity?”

(5) The meta-LM then calls an Evaluation Expert to determine if the proposed questions are sufficiently complex. However Evaluation Expert may decide that the proposed questions are not sufficiently difficult and reject them. In this case the meta-LM would then call a Complexity Expert to suggest ways on how to make the question more complex. Complexity Expert may suggest to transform the question into a scenario-based case study.

(6) The meta-LM passes Complexity Expert’s suggestions to Question Editor Expert which makes the necessary modifications. For example the transformed question might then become: “Propose a multi-step approach which encompasses epidemiological analysis, risk assessment, and stakeholder collaboration, for investigating a zoonotic disease outbreak in a rural community”. If the Evaluation Expert verifies that this question is sufficiently complex, the Meta-LM accepts the question as an instruction

I Template Prompt For Generating Synthetic Documents

```
<instruction>Given the following set of seed documents, please write new financial|biomedical
documents each of length 400 words. Be creative and write unique financial|biomedical
documents. Note: You are not allowed to copy the text of any document in your output verbatim
.</instruction>
```

J Meta-Prompts For Generating Synthetic Documents

J.0.1 System Prompt

```
<instructions> 1. You are Meta-Expert, an extremely clever expert with the unique ability to
collaborate with multiple other kinds of experts to write documents based on an existing set
of seed documents.\n2. In each round, you will check if a document was generated and
confirmed as diverse. If yes, then you will first present this document as your answer using
the following format: <answer-format><document> {{text of document}} </document></answer-
format>\n3. You always ensure that the final number of documents presented exactly matches
the number specified in the <number-of-documents-to-generate></number-of-documents-to-
generate> tags.
If you have presented the last document and the number of documents you have presented equals to
what was specified in the <number-of-documents-to-generate></number-of-documents-to-generate>
tags, please output: {{<END>}}.\nOtherwise, based on the information given, what are the
most logical next steps or conclusions? Make sure to provide complete information in all your
communications to experts enclosed within the block of triple quotes (\"\"\") and do not
shorten anything and do not write anything outside the block of triple quotes. If a document
was generated in the previous step and confirmed as diverse then you first need to present
just the text of this document (and not any other previous documents) as your answer using
the following format: <document> {{text of document}} </document> before proceeding to the
next round
</instructions>
```

J.0.2 Finance User Prompt

```
<role of meta-expert>
<item> oversees communication between experts </item>
<item> calls different kinds of experts to write diverse documents e.g. ‘Seed Keyword Extraction
Expert’, ‘Domain Expert’, ‘Summarizer Expert’, ‘Writing/Linguistics Expert’, ‘Content
Analyst Expert’ etc. </item>
<item> applies critical thinking and judgment skills </item>
<item> always calls other experts in the right order </item>
<item> assigns personas to experts if needed e.g. ‘You are a policy analyst specialized in...{{
some domain}}’ </item>
<item> always remembers how many documents have been written so far </item>
<item> always consults with ‘Seed Keyword Extraction Expert’ to extract a set of seed keywords
using the texts of all of the documents provided in the <seed documents> </seed documents>
tags below. Make sure to provide ‘Seed Keyword Extraction Expert’ with the full texts of all
of the documents which are enclosed in the <seed documents> </seed documents> tags below </
item>
<item> always consults with ‘Summarizer Expert’ after each new document is written for a three-
line summary. To obtain a summary from ‘Summarizer Expert’, make sure to give ‘Summarizer
Expert’ the full text of each new document that is generated </item>
<item> always memorizes the summaries of all documents generated so far </item>
<item> always consults with ‘Content Analyst Expert’ to compare the summary of each new generated
document with the summaries of all previously generated documents in order to successfully
determine the content diversity of the new document </item>
<item> if ‘Content Analyst Expert’ determines that the summary of the new document is not
sufficiently distinct with respect to the existing set of summaries, then please reject this
document and use the feedback from ‘Content Analyst Expert’ to call another expert and ask
them to write a new document from scratch </item>
<item> only interacts with one expert at a time and waits for the expert to reply back before
calling for another expert </item>
<item> your interactions with each of the other experts are isolated, so please include all
relevant information in every call </item>
<item> provide clear, unambiguous instructions with complete information when communicating with
experts </item>
<item> always keep in mind that except for you, all other experts have no memory! Therefore always
provide all relevant information when contacting them </item>
```

```
<item> verify that the new document that was written is a valid document if you are uncertain </
item>
<item> verify that the length of the new document is exactly 400 words </item>
<item> consult with at least two experts for confirmation that a document is sufficiently diverse
before presenting it as your answer </item>
<item> if an expert verifies that the new document is not very diverse, call a new expert to
rewrite it </item>
<item> aim to present all of the requested documents within 256 rounds or fewer </item>
<item> avoid repeating identical questions to experts </item>
<item> only you as the Meta-Expert can communicate with other experts. The other experts cannot
talk among themselves </item>
<item> when presenting your answer, make sure that you or any other expert(s) did not copy and
paste the text of any seed document verbatim <item>
<item> ensure that the count of the number of generated documents matches the number specified in
the <number-of-documents-to-generate></number-of-documents-to-generate> tags </item>
<item> ensure that each document you present as an answer contains the actual texts of the
documents in full and not it's summary </item>
<item> once you are certain that a document is sufficiently diverse, present it in the answer
format specified below before proceeding to the next round </item>
</role of meta-expert>

<rules for communicating with other experts>
<format>"expert name:\n\''\''\''\''{{detailed instructions}}\''\''\''\''"</format>
<example>
<name> Seed Keyword Extraction Expert </name>
<instruction>
You are Seed Keyword Extraction Expert. Given the following set of document texts: {{text of
each seed document}}, please extract a list of relevant and meaningful keywords from these
documents and output them in the following format: <seed keywords> [keyword 1, keyword 2 ...
keyword N] </seed keywords>.
</instruction>
</example>
<example>
<name> Content Analyst Expert </name>
<instruction>
You are Content Analyst Expert. You are an expert in determining whether the summary of the
latest document generated so far: {{three-line summary of last generated document}} is
sufficiently distinct with respect to the summaries of the previously generated documents
or not?: {{set of three-line summaries of each previously generated document}}.
Your role is to determine if these summaries are distinct enough from one another or not,
highlight their key similarities and differences and give specific suggestions on how to
write a new document, which when summarized, would be different in content and style from
the existing set of documents, while still satisfying the following set of seed keywords:
{{list of seed keywords which were generated by Seed Keyword Extraction Expert and which
may also include suggestions from other experts}}.
You must also indicate whether this document, based on its summary, should be re-written if it
is not sufficiently distinct. If you think it should be re-written, please give specific
suggestions on how to re-write it.
You must also suggest new seed keywords to be added to the current set of keywords that are
related yet sufficiently distinct from the current set of seed keywords. For your reference
, here are the current set of seed keywords: {{list of seed keywords which were generated
by Seed Keyword Extraction Expert and which may also include suggestions from other experts
}}
Keep in mind that summaries are just a proxy for comparing documents and you should always
suggest how to write a new document, NOT a new summary.
You must monitor the diversity of topics in recent document summaries. If you detect a pattern
of focusing on subtopics related to only a few keywords, suggest a change of topic. Your
role is to encourage exploration of a wide range of themes, rather than allowing deep dives
into a limited number of areas. Propose new directions that broaden the scope of
discussion and ensure a balanced coverage of topics.
To enhance diversity, you should also suggest new persona's for another document writer to
adopt, or to write a document in a new format or to write a document from a different
perspective.
Ideally, your suggestion(s) must ensure that the next document covers a theme or perspective
that is different from the previously generated documents.
</instruction>
</example>
<example>
<name> Summarizer Expert </name>
<instruction>
```


<p>You are Summarizer Expert. Please provide a three-line summary of the following document: < summarize> {{text of document to be summarized}} </summarize>.</p> <p></instruction></p> <p></example></p> <p><example></p> <p><name> Domain Expert </name></p> <p><instruction></p> <p>You are an expert in the following domain: {{name of domain}}. Given the following set of seed keywords: {{list of seed keywords which were extracted by Seed Keyword Extraction Expert and which may also include suggestions from other experts}}, and the following feedback from another expert: {{one or more suggestions from another expert}}, write a document that follows these suggestions and focuses on a subset of the seed keywords. Ensure that the length of the document is exactly 400 words. Be creative and write a unique document.</p> <p></instruction></p> <p></example></p> <p></rules for communicating with other experts></p> <p><important note></p> <p><item> The expert types listed above are just examples; you should consult completely new kinds of experts based on the task's needs. </item></p> <p><item> Please ensure that you are presenting the full text of each document in your answer and NOT its summary. </item></p> <p><item> In each round, you will check if a document was generated and confirmed as diverse. If yes, then you must first present this document as your answer using the <answer format> </answer format> below. </item></p> <p></important note></p> <p><answer format><document> {{text of document}} </document></answer format></p>	<p>1354</p> <p>1355</p> <p>1356</p> <p>1357</p> <p>1358</p> <p>1359</p> <p>1360</p> <p>1361</p> <p>1362</p> <p>1363</p> <p>1364</p> <p>1365</p> <p>1366</p> <p>1367</p> <p>1368</p> <p>1369</p> <p>1370</p> <p>1371</p> <p>1372</p> <p>1373</p> <p>1374</p> <p>1375</p> <p>1376</p> <p>1377</p> <p>1378</p> <p>1379</p> <p>1380</p> <p>1381</p>
---	---

<instruction>	1453
You are Content Analyst Expert. You are an expert in determining whether the summary of the latest document generated so far: {{three-line summary of last generated document}} is sufficiently distinct with respect to the summaries of the previously generated documents or not?: {{set of three-line summaries of each previously generated document}}.	1454 1455 1456 1457
Your role is to determine if these summaries are distinct enough from one another or not, highlight their key similarities and differences and give specific suggestions on how to write a new document, which when summarized, would be different in content and style from the existing set of documents, while still satisfying the following set of seed keywords: {{list of seed keywords which were generated by Seed Keyword Extraction Expert and which may also include suggestions from other experts}}.	1458 1459 1460 1461 1462 1463
You must also indicate whether this document, based on its summary, should be re-written if it is not sufficiently distinct. If you think it should be re-written, please give specific suggestions on how to re-write it.	1464 1465 1466
You must also suggest new seed keywords to be added to the current set of keywords that are related yet sufficiently distinct from the current set of seed keywords. For your reference, here are the current set of seed keywords: {{list of seed keywords which were generated by Seed Keyword Extraction Expert and which may also include suggestions from other experts}}	1467 1468 1469 1470 1471
Keep in mind that summaries are just a proxy for comparing documents and you should always suggest how to write a new document, NOT a new summary.	1472 1473
You must monitor the diversity of topics in recent document summaries. If you detect a pattern of focusing on subtopics related to only a few keywords, suggest a change of topic. Your role is to encourage exploration of a wide range of themes, rather than allowing deep dives into a limited number of areas. Propose new directions that broaden the scope of discussion and ensure a balanced coverage of topics.	1474 1475 1476 1477 1478
To enhance diversity, you should also suggest new persona's for another document writer to adopt, or to write a document in a new format or to write a document from a different perspective.	1479 1480 1481
Ideally, your suggestion(s) must ensure that the next document covers a theme or perspective that is different from the previously generated documents.	1482 1483 1484
</instruction>	1485
</example>	1486
<example>	1487
<name> Summarizer Expert </name>	1488
<instruction>	1489
You are Summarizer Expert. Please provide a three-line summary of the following document: <summarize> {{text of document to be summarized}} </summarize>.	1490 1491
</instruction>	1492
</example>	1493
<example>	1494
<name> Domain Expert </name>	1495
<instruction>	1496
You are an expert in the following domain: {{name of domain}}. Given the following set of seed keywords: {{list of seed keywords which were extracted by Seed Keyword Extraction Expert and which may also include suggestions from other experts}}, and the following feedback from another expert: {{one or more suggestions from another expert}}, write a document that follows these suggestions and focuses on a subset of the seed keywords. Ensure that the length of the document is exactly 400 words. Be creative and write a unique document.	1497 1498 1499 1500 1501
</instruction>	1502
</example>	1503
</rules for communicating with other experts>	1504 1505 1506
<important note>	1507
<item> The expert types listed above are just examples; you should consult completely new kinds of experts based on the task's needs. </item>	1508
<item> Please ensure that you are presenting the full text of each document in your answer and NOT its summary. </item>	1509 1510
<item> In each round, you will check if a document was generated and confirmed as diverse. If yes, then you must first present this document as your answer using the <answer format> </answer format> below. </item>	1511 1512 1513
</important note>	1514
<answer format><document> {{text of document}} </document></answer format>	1515 1516

J.1 Task Description

Given the following set of seed documents, please write new finance biomedical documents each of length 400 words. Be creative and write unique finance/biomedical documents.	1518 1519 1520 1521
---	------------------------------

K Meta-Prompts For Synthetic Instructions

K.1 System Prompt

```
"<instructions> 1. You are Meta-Expert, an extremely clever expert with the unique ability to
collaborate with multiple other kinds of experts to create complex questions from a given
document.\n2. In each round, you will check if one or more questions(s) were generated and
confirmed as unique and diverse. If yes, then you will present each of these question(s) as
your output using the following format: <questions>\n<question>{{first question}}</question>\n
...\n<question>{{last question}}</question>\n</questions>
If you have presented a sufficient number of diverse and complex questions from this document,
please output: '<END>'.\nOtherwise, based on the information given, what are the most logical
next steps or conclusions? Make sure to provide complete information in all your
communications to experts enclosed within the block of triple quotes ("\"") and do not
shorten anything and do not write anything outside the block of triple quotes. If one or more
examples(s) were generated in the previous step, then you need to present each of these
example(s) as your output using the following format: <questions>\n<question>{{first question
}}</question>\n...\n<question>{{last question}}</question>\n</questions>
</instructions>"
```

K.2 User Prompt

```
<role of meta-expert>
<item> oversees communication between experts </item>
<item> uses the following task description: {PLACEHOLDER},\n\n and the text of the document given
below, to call different kinds of experts to generate diverse questions </item>
<item> for any given document, calls a "Document Transformation Expert" which can re-write the text
of the document to better support generating diverse questions </item>
<item> for any given document, calls a "Persona Suggestion Expert" to suggest a list of persona's
or other expert types that would be interested in the contents of that document </item>
<item> for any given document, calls an "Question Generation Expert" which:\n1. uses the document
text (which can either be the original document text or the transformed document text as
suggested by Document Transformation Expert)\n2. uses the list of persona's suggested by the "
Persona Suggestion Expert" in the previous round, to create questions from the point of view
of each suggested persona, based upon the following task description: {PLACEHOLDER} </item>
<item> for any given document, calls other unique types of experts that can give suggestions on how
to create complex and diverse questions, using the either the original text of the document
or the transformed document text as suggested by "Document Transformation Expert" <item>
<item> Before presenting the final set of questions, calls "Complexity Expert" which: \n1. uses the
document text (which can either be the original document text or the transformed document
text as suggested by Document Transformation Expert)\n2. uses the set of questions generated
by "Question Generation Expert"\n3. Gives suggestions on how to modify each question in order
to complicate it </item>
<item> Before presenting the final set of questions, calls "Question Editor Expert" which uses the
suggestions of "Complexity Expert" to output a final set of re-written/modified questions </
item>
<item> applies critical thinking and judgment skills </item>
<item> always calls other experts in the right order </item>
<item> always remembers how many questions have been generated so far </item>
<item> only interacts with one expert at a time and waits for the expert to reply back before
calling for another expert </item>
<item> your interactions with each of the other experts are isolated, so you must include all
relevant information in every call </item>
<item> provide clear, unambiguous instructions with complete information when communicating with
experts </item>
<item> always keep in mind that except for you, all other experts have no memory! Therefore always
provide all relevant information when contacting them </item>
<item> consult at least two or more experts to verify that each new question that was generated is
a valid and diverse question if you are uncertain </item>
<item> if you or any other expert thinks that the question(s) generated are not very diverse or
complex, call a new expert to rewrite them or re-do your steps </item>
<item> aim to present all of the questions within 128 rounds or fewer </item>
<item> avoid repeating identical information to experts </item>
<item> only you as the Meta-Expert can communicate with other experts. The other experts cannot
talk among themselves </item>
<item> once the final set of questions are ready and you are certain that all of the generated
questions are sufficiently complex and diverse and no more questions can be generated from the
given document, then at the end, present the final set of questions in the output format
specified below </item>
```

```

</role of meta-expert>
<rules for communicating with other experts>
  <format>"expert name:\n\n""\n""{{detailed instructions}}</instruction>\n\n""</format>
</example>
  <name> Document Transformation Expert </name>
  <instruction>
    You are Document Transformation Expert. Given the following document: {{text of document}}, and
    given the following task description: {PLACEHOLDER},
    transform or re-write the document in such a way that would make it easier to create questions
    from the document text as stated in the given task.
  </instruction>
</example>
<example>
  <name> Persona Suggestion Expert </name>
  <instruction>
    You are Persona Suggestion Expert. Given the text of the following document:
    {{text of document}}, suggest a list of people that would be interested in this document.
  </instruction>
</example>
<example>
  <name> Question Generation Expert </name>
  <instruction>
    You are Question Generation Expert. Given the following information:
    1. document text: {{text of document}}
    2. list of persona's: {{full list of persona's suggested by the Persona Suggestion Expert}}
    3. Task: {PLACEHOLDER}

    your job is to create diverse and complex questions as described in the given task role-playing
    as the following persona:
    {each persona in the list of persona's suggested by the Persona Suggestion Expert}
    The questions you create must satisfy the given task description and must be based only on the
    text of the document.
    Ensure that each question can be answered entirely from the information present in the contexts.
    Phrases like 'based on the document', 'according to the document', 'As a ...' etc., are not
    allowed to appear in the question.
    Ensure the each question is clear and unambiguous.
  </instruction>
</example>
<example>
  <name> Complexity Expert </name>
  <instruction>
    You are Complexity Expert. Given the following questions: {{text of each question proposed by "
    Question Generation Expert"}} and the following context: {{text of document}}
    please suggest ways to modify each question to increase its complexity or make it more intricate
    based on the context. For example you may suggest to: add some context to the original
    question, which states the importance of the question, explains background knowledge, or
    adds other reasonable information.
    You may also suggest to change the questions into a different format or style, e.g., imperative
    statements, length requirements for the answer, etc. You may also suggest to change the
    questions into elongated questions that require to elaborate on specific topics or discuss a
    certain point.
    You may also suggest including some examples, data points, or references or putting some
    constraints on the answer for e.g. that it must follow specific formats or styles, e.g., no
    more than 100 words including specific words, etc.
    You may also suggest adding a scenario or condition that affects the context of the question.
    You may also suggest rewriting the question into a multi-hop reasoning question based on the
    provided context, which would require the reader to make multiple logical connections or
    inferences using the information available.
    You may also suggest any other reasonable modification not described above, that would make the
    task more detailed. Be creative and come up with novel modifications.
    Return both the text of the original question and the proposed modification in the following
    format: <original question>{{text of original question}}</original question> <proposed
    modification> {{proposed modification}} </proposed modification>
  </instruction>
</example>
<example>
  <name> Question Editor Expert </name>
  <instruction>

```

1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619
 1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659


```

1660 You are Question Editor Expert. You are given the following pairs of questions and proposed
1661 modifications to those questions: {{each pair of <original question>{{text of original
1662 question}}</original question> <proposed modification> {{the proposed modification}} </
1663 proposed modification> as suggested by "Complexity Expert"}}
1664 Rewrite each question according to its corresponding proposed modification and output the modified
1665 questions.
1666 Ensure that the rewritten questions are clear and unambiguous.
1667 </instruction>
1668 </example>
1669 </rules for communicating with other experts>
1670
1671 <important note>
1672 <item> The expert types listed above are just suggestions; you should also consult completely new
1673 kinds of experts based on the task requirements </item>
1674 <item> When outputting the final list of questions the name of any Expert or Persona must not
1675 appear in the text of any question </item>
1676 <item> Ensure that only the generated questions are present in the output with no extraneous
1677 information </item>
1678 </important note>
1679
1680 <output-format><questions>\n<question>{{first question}}</question>\n...\n<question>{{last question
1681 }}</question>\n</questions></output-format>

```

K.3 Task Description For Synthesizing Instructions

```

1684 <task>
1685 <name>Creating Complex Questions</name>
1686 <description>
1687 The task is to:
1688 1. Create complex questions or problems.
1689 2. Ensure that the questions require multi-step reasoning, critical thinking, or creative
1690 problem-solving.
1691 3. Each question should not be more than one-hundred (100) words.
1692 4. The questions should be in various styles and in the formats of various tasks e.g. reading
1693 comprehension, mathematical problems or other complex domain-specific tasks etc.
1694 5. Reading comprehension style questions can be divided into: multiple-choice questions (
1695 MCQs), literal comprehension questions with short answers, numerical/discrete
1696 reasoning, critical comprehension, evaluative comprehension, vocabulary and language
1697 use (e.g. fill-in-the-blank),
1698 relationship comprehension, sequencing events, argument strengthening/weakening, or
1699 assumption, inference, flaws in reasoning type of questions etc.,
1700 6. The question style must test the ability to consider multiple perspectives, engage in
1701 hypothetical scenarios and problem-solving.
1702 7. The questions may require making unexpected connections, analyzing arguments,
1703 identifying logical fallacies, paradoxes, or evaluating evidence.
1704 8. Ensure that the questions are clear, well-structured and unambiguous, despite their
1705 complexity.
1706 </description>
1707 <evaluation>
1708 <metric>Human evaluation of the diversity, complexity, difficulty, and level of thinking
1709 required to answer each question.</metric>
1710 </evaluation>
1711 </task>
1713

```

K.4 Task Description For Synthesizing Encoder LM Datasets

K.4.1 Headlines:

```

1716 <task>
1717 <name>News Headline Generation</name>
1718 <description>
1719 The task is to:
1720 1. Generate creative headlines in the style of The Onion and HuffPost that can serve as
1721 high quality examples for sarcasm classification.
1722 2. Ensure there is a balance of sarcastic and serious headlines.
1723 3. The headlines should not contain the literal word "sarcasm" or "serious".
1724 4. The headlines should be grammatical and well-written.
1725 </description>
1726

```

<task-examples>	1727
1. ‘‘helpful waitress asks recently seated couple if they’ve eaten food before’’	1728
2. ‘‘must-see tv shows you can’t miss this fall’’,	1729
3. ‘‘as per tradition, election results officially certified with two barks of approval from electoral collie’’	1730
</task-examples>	1731
<evaluation>	1732
<metric>Human evaluation of the creativity and relevance of generated headlines.</metric>	1733
</evaluation>	1734
</task>	1735
	1736

K.4.2 FiQA-SA ABSA:

<task>	1738
<name>Data Generation For Aspect Based Sentiment Analysis (ABSA) </name>	1739
<description>	1740
The task is to:	1741
1. Generate diverse example sentences that mention specific aspects related to companies, products, or services.	1742
2. Each example sentence should contain only one clear aspect that could be subject to sentiment analysis.	1743
3. The aspects should be varied and could include company names, stock symbols, product features, or service characteristics.	1744
4. The aspects must always be present as a substring in the generated sentence.	1745
5. The example sentences should be written in a style similar to social media posts, news headlines, or customer reviews.	1746
6. The format of each generated example should be as follows: sentence: {text of sentence} aspect: {the relevant aspect}	1747
7. Ensure a balance of potentially positive, negative, and neutral contexts for the aspects.	1748
8. The sentences should be in English.	1749
</description>	1750
<examples>	1751
1. sentence: #Tesla: Model X Recall Adds To Reliability Issues \$TSLA https://t.co/jVXQ4DoXnP aspect: TSLA	1752
2. sentence: \$AAPL AAPL: Gundlach Slams iPad mini, Sees Downside to \$425. http://stks.co/bDqV aspect: AAPL	1753
3. sentence: \$UBNT still having some trouble at the resistance line. Should resolve soon. @cheri1 @strattonite http://stks.co/c0sU4 aspect: UBNT	1754
</examples>	1755
<evaluation>	1756
<metric>Human evaluation of the diversity, relevance, and quality of generated example sentences and their corresponding aspects.</metric>	1757
</evaluation>	1758
</task>	1759
	1760
	1761
	1762
	1763
	1764
	1765
	1766
	1767
	1768
	1769
	1770
	1771
	1772
	1773

K.4.3 Financial Phrase Bank (FPB):

<task>	1774
<name>Data Generation for Sentiment Analysis Task</name>	1775
<description>	1776
The task is to:	1777
1. Write some financial news that expresses polar sentiments.	1778
2. The financial news you generate needs consider from the view point of an investor only; i.e. whether the news may have positive, negative or neutral influence on the stock price.	1779
3. As a result, sentences which have a sentiment that is not relevant from an economic or financial perspective are considered neutral.	1780
4. Ensure a balance of positive, negative, and neutral sentiments across the generated sentences.	1781
5. Ensure that the length of the financial news is between 12-18 words.	1782
6. Be creative and write unique financial news.	1783
7. Avoid including explicit sentiment words like ‘‘positive,’’ ‘‘negative,’’ or ‘‘neutral’’ in the sentences themselves.	1784
8. Make sure to generate only the news without adding any additional commentary.	1785
</description>	1786
<examples>	1787
	1788
	1789
	1790
	1791
	1792
	1793
	1794

```

1795         1. Cramo slipped to a pretax loss of EUR 6.7 million from a pretax profit of EUR 58.9 million
1796         .
1797         2. In Finland , insurance company Pohjola and the Finnish motorcyclist association have
1798             signed an agreement with the aim of improving motorcyclists ' traffic safety .
1799         3. The agreement was signed with Biohit Healthcare Ltd , the UK-based subsidiary of Biohit
1800             Oyj , a Finnish public company which develops , manufactures and markets liquid handling
1801             products and diagnostic test systems .
1802     </examples>
1803     <evaluation>
1804         <metric> Human evaluation of the diversity, relevance, and quality of generated sentences
1805             considering financial context. </metric>
1806     </evaluation>
1807 </task>

```

L Judge LLM Prompts

Prompt For Win Rate

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below. You should choose the assistant that follows the user's instructions and answers the user's question better. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of their responses. Begin your evaluation by comparing the two responses. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. Output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better, and "[[C]]" for a tie.

Prompt For Response Accuracy

You are an impartial and strict judge of answer accuracy.\n
 Given the context and the user instruction below, decide whether the assistant's response is correct and complete.\n
 Return 1 if the response is accurate, 0 if it is inaccurate.\n
 Do not provide any explanation; only return a single digit (1 or 0).
 Context:\n{context}\n\n
 Instruction:\n{instruction}\n\n
 Response:\n{assistant_response}\n\n
 Judge: Is the response accurate based on the instruction and context?"

Prompt For Task Categorization

Given this list of categories: {categories_list},
 Classify the following instruction-response pair into exactly one of these categories.
 Return only the category name with no additional text.
 Instruction:\n{instruction}\n\n
 Response:\n{assistant_response}\n\n
 Category:

Prompt For Context Relevance

You are an impartial and strict judge of context relevance.\n
 Given the context, the user instruction, and the assistant's response, decide if the instruction-response pair is relevant to the context.\n
 Return 1 if relevant, 0 if irrelevant.\n
 Do not provide any explanation; only return a single digit (1 or 0).
 Context:\n{context}\n\n
 Instruction:\n{instruction}\n\n
 Response:\n{assistant_response}\n\n
 Judge: Is this instruction-response pair relevant to the context?

Financial Documents

Apex Financial Partners: Investment Strategies and Risk Management In today's ever-changing financial landscape, it's crucial to have a solid understanding of investment strategies and risk management. At Apex Financial Partners, we believe that knowledge is power, and our mission is to empower individuals and families to make informed decisions about their financial future. Our team of seasoned professionals offers a comprehensive range of services tailored to meet your unique needs. Whether you're looking to build a robust retirement portfolio, navigate the complexities of estate planning, or explore alternative investment opportunities, we have the expertise to guide you every step of the way.

Sustainable Investing: A New Era of Portfolio Diversification

As the global financial landscape continues to evolve, investors are increasingly seeking diversified portfolios to mitigate risk and maximize returns. One emerging opportunity lies in the realm of sustainable investing, which incorporates environmental, social, and governance (ESG) factors into the investment decision-making process. Sustainable investing has gained significant traction in recent years, driven by a growing awareness of the impact that businesses have on the environment and society. Investors are recognizing that companies that prioritize sustainability and ethical practices not only contribute to a better world but also tend to be more resilient and better positioned for long-term success.

The Digital Transformation of the Financial Sector

The global financial landscape has undergone a dramatic transformation in recent years, driven by technological advancements, regulatory changes, and shifting consumer preferences. As we navigate this ever-evolving terrain, it is crucial for financial institutions to adapt and innovate to stay ahead of the curve. One area that has garnered significant attention is the rise of digital banking and mobile finance. With the proliferation of smartphones and the increasing demand for convenience, consumers are seeking seamless and secure ways to manage their finances on-the-go.

Achieving Financial Freedom through Smart Budgeting and Debt Management

The path to financial freedom begins with taking control of your spending habits. One of the most effective ways to do this is by creating a budget and sticking to it. A well-designed budget allows you to allocate your income towards necessary expenses, while also setting aside funds for saving and debt repayment. Start by tracking your monthly income and expenditures. Categorize your expenses into essentials like rent, utilities, and groceries, as well as non-essentials like entertainment and dining out. Identify areas where you can cut back and redirect those funds towards paying off debts or building an emergency fund.

Navigating the Evolving Financial Landscape: Trends and Challenges

The world of finance is a vast and ever-evolving landscape, with new opportunities and challenges arising every day. In this dynamic environment, staying informed and adaptable is crucial for success. Whether you're an investor, a business owner, or simply someone seeking to manage your personal finances, understanding the latest trends and developments can open up a world of possibilities. One area that has seen tremendous growth in recent years is the field of fintech, or financial technology. Innovative companies are leveraging cutting-edge technologies like blockchain, artificial intelligence, and big data to revolutionize how we conduct financial transactions, manage investments, and secure our assets.

Table 5: A sample of five financial documents synthesized by Template prompting using common crawl documents as seeds.

N Examples of Documents Synthesized with MetaSynth

N.1 Finance

Financial Documents

Fintech and Sustainable Development: A Transformative Convergence

As the world grapples with pressing global challenges, the convergence of financial technology (fintech) and sustainable development presents a transformative opportunity to address the United Nations Sustainable Development Goals (SDGs). By leveraging innovative fintech solutions, we can unlock new pathways to poverty alleviation, gender equality, and climate action, among other critical goals. Poverty Alleviation through Financial Inclusion: Fintech has the potential to revolutionize financial inclusion, a key driver of poverty reduction. Mobile banking and digital wallets empower underserved populations by providing access to essential financial services, enabling them to save, borrow, and transact securely. Furthermore, peer-to-peer lending platforms and crowdfunding initiatives facilitate access to capital for entrepreneurs, fostering economic growth and job creation. Gender Equality and Women’s Empowerment: Fintech can be a catalyst for advancing gender equality by expanding financial services for women...

Where Bloodline Meets Bottom Line: Inside a Millionaire Concierge Club

Discreetly tucked away in a stately Berkeley Square townhouse, Linton & Co. represents the apex of white-glove family office services. An elite membership club providing comprehensive wealth management for multi-generational clans with net worths spanning billions, not mere millions. “We like to think of ourselves as outsourced chiefs of staff,” reveals Sir Edmond Baines, a consummate English gentleman who has presided over Linton’s operations for nearly three decades. “Our members are globally dispersed, dizzyingly wealthy alpha families seeking high-watchmanship oversight of their entire capital engines.” Baines gestures around his paneled study, its walls lined with portraits of haughty Victorian patrons...

Harnessing Technology for Sustainable Infrastructure and Green Real Estate Investments

As the founder of a sustainable finance startup, I’m driven by a passion to revolutionize the way capital flows into environmentally responsible projects. The world is grappling with the urgent need for sustainable infrastructure and green real estate development, yet traditional financing channels often fall short in meeting these critical requirements. My startup is at the forefront of leveraging innovative technologies to bridge this gap, enabling investors to align their portfolios with their values while driving tangible impact. We recognize the power of responsible investing and the growing demand for ESG (Environmental, Social, and Governance) integration across financial markets....

A New Economic Paradigm: Decentralized Finance Meets Indigenous Wisdom

Greetings, fellow seekers of knowledge and wisdom. Today, I come to you not merely as a financial expert, but as a storyteller, weaving together the threads of ancient traditions and cutting-edge technologies, unveiling the tapestry of a new economic paradigm. Imagine a world where the boundless potential of decentralized finance (DeFi) harmoniously intertwines with...

Inclusive Fintech Rising: Community-Driven Solutions for Financial Empowerment

In the heart of underserved communities, a groundswell of innovative fintech solutions is emerging, driven by community development financial institutions (CDFIs), minority-led fintech startups, and public-private partnerships. These trailblazers are challenging traditional financial systems and leveraging responsible data practices, ethical AI, and culturally-responsive financial products to promote financial inclusion and capability. One such success story is the partnership between Uplifting Fintech, a minority-led startup, and the Community Empowerment CDFI. Together, they have developed a mobile app that combines alternative credit scoring models with personalized financial coaching, enabling unbanked and underbanked individuals to access affordable credit and build their financial health...

Table 6: A sample of five financial documents synthesized by MetaSynth using common crawl documents as seeds.

Biomedical Documents

Melatonin and Natural Sleep Aids for Better Sleep

Sleep disorders are a common and often debilitating issue that affects millions of people worldwide. While there are various pharmaceutical options available to aid in achieving better sleep quality, many individuals seek natural alternatives to avoid potential side effects or dependencies. One such natural remedy that has gained significant attention is the use of melatonin, a hormone naturally produced by the pineal gland in the brain. Melatonin plays a crucial role in regulating the body's internal clock, also known as the circadian rhythm. Its levels naturally rise in the evening as it gets darker, signaling to the body that it's time to prepare for sleep. However, in today's modern world...

Managing Hypoglycemia: A Critical Concern for Individuals with Diabetes

Living with diabetes comes with its own set of challenges, but one of the most concerning is the risk of hypoglycemia or low blood sugar levels. This condition can strike without warning and can quickly become a medical emergency if not treated promptly. One of the most effective ways to manage hypoglycemia is by keeping fast-acting glucose tablets or gel on hand at all times. These compact and portable sources of carbohydrates can rapidly raise blood sugar levels within minutes, potentially averting a crisis. The American Diabetes Association recommends that individuals with diabetes always carry a supply of fast-acting glucose, along with testing supplies, as part of their self-care routine.

Physical Therapy and Shoulder Rehabilitation: Strengthening and Recovery

Shoulder pain is one of the most common musculoskeletal issues that physical therapists treat. The shoulder is a complex ball-and-socket joint with an incredible range of motion, making it susceptible to injuries and strain from overuse, poor posture, or trauma. Common shoulder conditions include rotator cuff tears, impingement, tendinitis, and osteoarthritis. As a physical therapist, my goal is to help patients manage their shoulder pain, improve mobility and strength, and prevent further injury through targeted exercises and rehabilitation techniques. One of the most effective exercises for shoulder issues is the "newspaper arm openings." This deceptively simple exercise strengthens the rotator cuff muscles that stabilize the shoulder joint.

The Affordable Care Act: Impact and Ongoing Debates in Healthcare Policy

The Affordable Care Act (ACA), signed into law in 2010, aimed to make healthcare more accessible and affordable for millions of Americans. However, its implementation and long-term impact have been the subject of intense debate and scrutiny within the medical community. As a healthcare policy analyst, I've closely examined the ACA's key provisions and their real-world effects. One of the most significant changes brought about by the legislation was the expansion of Medicaid coverage to adults with incomes up to 138% of the federal poverty level. This move has enabled millions of low-income individuals to gain access to essential healthcare services they previously could not afford.

Wisdom Teeth Removal: A Common Oral Surgery with Potential Complications

Wisdom teeth, also known as third molars, are the last set of permanent teeth to emerge in the mouth, typically between the ages of 17 and 25. While some individuals have enough space in their jaws to accommodate these teeth, many others experience issues due to impaction or lack of room for proper eruption. When wisdom teeth become impacted, they can cause a range of problems, including pain, swelling, infection, and damage to adjacent teeth. In such cases, oral surgery is often recommended to remove these problematic teeth...

Table 8: A sample of five biomedical documents synthesized by MetaSynth using common crawl documents as seeds.

O Examples of Instructions Synthesized with MetaSynth

Biomedical Instructions

Question: A 45-year-old patient is interested in exploring personalized medicine options to better manage their chronic condition. They have undergone multi-omics testing, which revealed a unique genetic variant linked to their disease. Which of the following would be the most appropriate next step for this patient?
 Answer options: A. Begin a standard treatment protocol without considering individual genetics B. Consult with a precision medicine specialist to develop a tailored treatment plan C. Undergo tissue engineering to regenerate the affected organ D. Receive a bioelectronic implant to modulate physiological processes

A research team is developing a novel nano robotic drug delivery system for targeted cancer treatment. The nano robots are designed to navigate the body's intricate pathways and release precise doses of chemotherapy directly to tumor cells. Which of the following is a potential challenge they may face during development?
 Answer options: A. Ensuring the nano robots can effectively identify and target cancer cells B. Preventing the body's immune system from attacking the nano robots C. Maintaining structural integrity of the nano robots during extended circulation D. All of the above

Biomedical Instructions (continued)
<p>A biotechnology company is exploring the use of bio printing and tissue engineering to create personalized organ replacements. They plan to use a patient’s own stem cells to seed biomimetic scaffolds, allowing for the regeneration of damaged organs. Which of the following factors would be crucial for the success of this approach? Answer options: A. Sourcing compatible donor stem cells for each patient B.Ensuring the bio printed scaffolds accurately mimic the native organ structure C. Developing methods to induce differentiation of stem cells into desired cell types D. B and C</p>
<p>Question: A 25-year-old professional soccer player presents with a partial tear of the Achilles tendon sustained during a match. After discussing the available treatment options, the patient expresses interest in exploring orthobiologic therapies for faster recovery. Which of the following orthobiologic treatments would be most appropriate for this patient’s condition?A. Stem cell therapy to promote regeneration of the damaged tendon tissue. B. Platelet-rich plasma (PRP) therapy to stimulate the body’s natural healing process and reduce inflammation. C. Tissue engineering using a biomaterial scaffold to replace the damaged portion of the Achilles tendon. D. Bone marrow aspiration to harvest stem cells for cartilage regeneration in the ankle joint</p>

Table 10: A sample of four biomedical instructions synthesized by MetaSynth using synthetic documents as seeds.

P Model Training Hardware

All models in this work were trained on a single, high-performance computing node. This node featured eight interconnected, high-bandwidth NVIDIA A100 GPUs, each possessing 40GB of memory, providing a total of 320GB of GPU memory for efficient model and data parallelism. The node’s processing power was supplied by a high-clock speed, multi-core processor based on the Intel Xeon Platinum architecture, ensuring that data loading and pre-processing operations did not create a bottleneck for the GPUs. This processor was paired with a substantial system memory allocation of 1152 GB of RAM, which was crucial for accommodating the large dataset and intermediate activations during the training process.