DATA-CENTRIC LESSONS TO IMPROVE SPEECH-LANGUAGE PRETRAINING

Anonymous authors

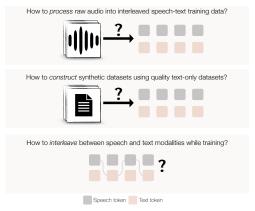
Paper under double-blind review

ABSTRACT

Spoken Question-Answering (SQA) is a core capability for useful and interactive artificial intelligence systems. Recently, several speech-language models (SpeechLMs) have been released with a specific focus on improving their SQA performance. However, a lack of controlled ablations of pretraining data processing and curation makes it challenging to understand what factors account for performance, despite substantial gains from similar studies in other data modalities. In this work, we address this gap by conducting a data-centric exploration for pretraining SpeechLMs. We focus on three questions fundamental to speech-language pretraining data: (1) how to *process* raw web-crawled audio content for speech-text pretraining, (2) how to *construct* synthetic datasets to augment web-crawled data and (3) how to *interleave* (text, audio) segments into training sequences. We apply the insights from our controlled data-centric ablations to pretrain a 3.8B-parameter SpeechLM, called **SpeLangy**, that outperforms models that are up to 3x larger by 10.2% absolute performance. We hope our findings highlight the impact of effective data curation and guide future data-centric exploration in SpeechLMs.

Our Three Data-Centric Research Questions

Our SpeLangy model is small and highly performant



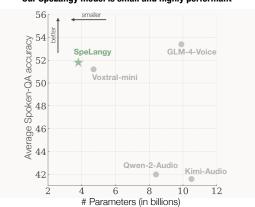


Figure 1: (**Left**) We highlight the *data-centric* questions we study in this work (Sec. 3), (**Right**) Distilling all our data-insights yields a strong 3.8B-parameter SpeechLM, **SpeLangy** (Sec. 5).

1 Introduction

Language-based assistants are now widely deployed (OpenAI, 2024; Comanici et al., 2025). Yet, purely textual interactions are inherently limiting for real-world assistants that must operate in open, hands-free settings. Voice provides a natural, low-friction interface for human—AI interaction, and recent work therefore emphasizes Spoken Question-Answering (SQA) (Nachmani et al., 2023; Liu et al., 2025; Xiaomi, 2025)—where a question is asked in audio and the system must produce spoken or textual answers—as a core capability for end-to-end speech language models (SpeechLMs).

Recently, *speech–text interleaved pretraining*—next-token prediction over sequences that alternate between speech and text tokens—has been proposed as a viable strategy to boost SQA performance (Nguyen et al., 2025b; Zeng et al., 2024b). However, while these works describe modeling choices

comprehensively, details of their data pipelines are often not shared or evaluated in a controlled setting. How should we process raw audio into trainable speech-text chunks? Can we leverage text-only datasets to go beyond datasets sourced from raw audio? How should we interleave tokens for effective modality alignment? In the current literature, *these data-centric questions remain underexplored*. In other domains like language (Dubey et al., 2024; Li et al., 2024) and vision (Gadre et al., 2023; Siméoni et al., 2025), data curation has consistently proven to be a primary driver of performance improvements, yet a large gap exists from the data-centric perspective in the speech-language domain.

In our work, we aim to close this gap with a systematic, *data-centric* study of interleaved pretraining for SQA (Fig. 1). We first provide a detailed description of our processing pipeline for converting raw audio into speech-text interleaved data (Fig. 8). We then study optimal interleaving strategies for speech-text pretraining, finding that fine-grained interleaving (which alternates between speech and text modalities at sentence boundaries) improves alignment of the two modalities (Sec. 3.3). Building on this, we introduce effective synthetic data methods involving LLM-based rewriting and text-to-speech synthesis to go beyond raw web-crawled audio for pretraining (Sec. 3.4). We also examine two modality-sampling schemes for interleaved training, finding that a deterministic ordering of alternating speech-text chunks is beneficial compared to stochastic modality sampling (Sec. 3.5). To understand *why* our data-centric methods improve SQA performance, we analyse the modality gap between speech and text distributions (Sec. 4.1) and inspect the topic distributions of web-crawled and synthetic datasets (Sec. 4.2). Finally, to showcase the efficacy of our data interventions at scale, we pretrain a 3.8B SpeechLM (SpeLangy) that outperforms 3x larger models by upto 10% average SQA performance. Taken together, our results underscore the central role of data curation in speech–language pretraining and motivate a broader, systematic push toward data-centric exploration.

2 RELATED WORK

Speech Language Models. Most recent SpeechLMs employ a simple Speech Encoder + Connector + LLM philosophy for conducting joint speech-text training (Lakhotia et al., 2021; Algayres et al., 2023; Hassid et al., 2023; Nguyen et al., 2025b; Nachmani et al., 2023; Rubenstein et al., 2023; Zhang et al., 2023; Défossez et al., 2024). Models like Kimi-Audio (Ding et al., 2025), Step-Audio-2 (Wu et al., 2025), Baichuan-Audio (Li et al., 2025b), GLM-4-Voice (Zeng et al., 2024a), and MiMo-Audio (Xiaomi, 2025) have emerged as strong foundation models that seamlessly perform several tasks, including spoken question-answering. While demonstrating impressive performance, details behind their data curation strategies are however scant. Through our controlled experiments, we aim to fill this gap by shedding light on how to effectively construct speech-text pretraining datasets.

Data Curation for Foundation Models. Pretraining data quality is pivotal for driving performance of foundation models. Efforts like Gopher (Rae et al., 2021), T5 (Raffel et al., 2020), Nemotron-CC (Su et al., 2024), FineWeb (Penedo et al., 2024), DCLM (Li et al., 2024) and Olmo-2 (OLMo et al., 2024) significantly emphasize the benefits of strong data processing, curation and filtering for language data. In computer vision, Dino-v2 (Oquab et al., 2023), Dino-v3 (Siméoni et al., 2025), Aim-v2 (Fini et al., 2025) and Web-SSL (Fan et al., 2025) showcased the high impact that careful data curation has on model quality. Similar results on the importance of data-centric research have been shown in vision-language (Gadre et al., 2023; Fang et al., 2023a; Tong et al., 2024) and reasoning-based (Guha et al., 2025; Li et al., 2025c; Muennighoff et al., 2025) foundation modeling literature. Owing to the paucity of such data-centric research in the speech-language domain, we aim to close this gap, demonstrating the strong utility of data-centric approaches for boosting SpeechLM quality.

3 CONTROLLED DATA-CENTRIC EXPERIMENTS

In this section, we address three *data-centric* questions for improving SQA, via controlled experiments: (1) how to *process* raw web-crawled audio into suitable interleaved speech-text training data (Sec. 3.3), (2) how to *construct* synthetic speech-text datasets seeded from text-only datasets (Sec. 3.4), and (3) how to *interleave* between speech and text modalities while training (Sec. 3.5).

3.1 EVALUATION BENCHMARKS

Spoken Question-Answering ($S \rightarrow T$). We use three standard benchmarks for SQA where the model is asked questions in speech and responds in text ($S \rightarrow T$): *Spoken-LLaMA-Questions* (SLQ),

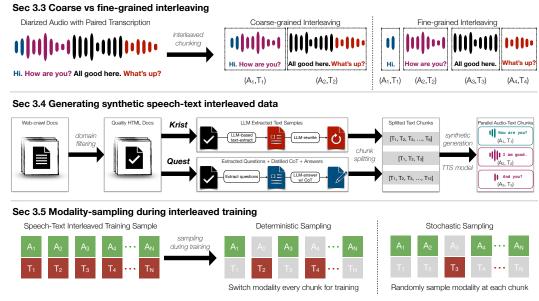


Figure 2: **Our experimental conditions for speech-text pretraining data.** (**Top**) We study two interleaving strategies: *coarse* (long chunks) and *fine* (short chunks) (Sec. 3.3). (**Middle**) We construct two synthetic datasets—*Krist* and *Quest*—from web-crawled documents (Sec. 3.4). (**Bottom**) We study two modality-sampling schemes: *deterministic* and *stochastic* (Sec. 3.5).

Spoken-Web-Questions (SWQ) and Spoken-TriviaQA (STQ). We source audio questions from OpenAudioBench (Li et al., 2025b). Our protocol follows standard language modeling pretraining evaluations (Gu et al., 2024; Allal et al., 2025) to use an MCQ cloze-format with log-likelihood evaluation for choosing the correct option (we use 4 multiple choices with chance-level accuracy being 25%). We provide more details and examples from each of our evaluation sets in Appx. F.

Text Understanding (T\rightarrowT). To ensure our speech-text pretraining recipe does not degrade base language performance, we evaluate on 12 standard text benchmarks spanning across general knowledge, math and coding: *MMLU* (Hendrycks et al., 2020), *CoreEN* (Gunter et al., 2024; Mizrahi et al., 2025) (consisting of 9 benchmarks—*ARC-Easy* and *ARC-Challenge* (Clark et al., 2018), *HellaSwag* (Zellers et al., 2019), *Lambada* (Paperno et al., 2016), *PIQA* (Bisk et al., 2020), *SciQ* (Welbl et al., 2017), *TriviaQA* (Joshi et al., 2017), *WebQuestions* (Berant et al., 2013), and *WinoGrande* (Sakaguchi et al., 2021)), *GSM-8k* (Cobbe et al., 2021), and HumanEval (Chen et al., 2021).

3.2 Base Setup

Model Architecture. We conduct all our experiments with a \sim 3.8B-parameter SpeechLM, consisting of two major components: a speech tokenizer and a pretrained language model. Our speech tokenizer consists of a speech encoder with conformer (Gulati et al., 2020) blocks followed by a finite scalar quantizer (Mentzer et al., 2023) that outputs discrete speech tokens. We initialize our LM with the dense 3B base-LM from (Li et al., 2025a) that has a context-length of 16, 384 tokens.

Training Data. Our base mixture consists of web-crawled audio that we process into interleaved speech-text data. We provide more details on how we process audio into our training data format in the next section. We also use the text continued-pretraining dataset from (Li et al., 2025a) to preserve the base-LM's text performance. Following prior multimodal works (Shukor et al., 2025; McKinzie et al., 2024), we use a 60% text-only and 40% speech-text data mixture during training.

Optimization Details. We train with a global-batch-size of 512 and a packed-sequence-length of 16, 384 tokens, for 200k steps. We use the standard next-token prediction objective and compute the loss over both speech and text tokens (we also conduct ablations with loss-masking on the speech tokens in Sec. 3.6). We only tune the language model weights while keep the speech tokenizer frozen. For more details, please refer to Appx. D.

163 164

165

166

167

168

170 171

172

173

174

175

176

177

178

179

181

182

183

185

187

196 197

199

200

201

202

203

204

205

206

207 208

209

210

211

212

213

214

215

3.3 Processing pretraining data via fine-grained interleaving

Extracting interleaved data from raw audio. We begin with >10M hours of raw web-crawled audio. To process them into trainable speech-text samples, we follow a multi-stage pipeline (see Fig. 8), involving speaker diarization, language detection and filtering, paired-transcription generation and filtering, and interleaved chunking. Our pipeline yields interleaved training samples X_i consisting of multiple paired speech-text chunks of the form $X_i = \{(A_1, T_1), (A_2, T_2) \cdots (A_n, T_n)\}$, where n is the number of chunks in each sample. We provide more details about each individual component along with detailed statistics in Appx. A, while focusing on the interleaved chunking component here.

Fine vs coarse interleaving. Prior speech-text pretraining works (Liu et al., 2025; Zeng et al., 2024a) have explored constructing interleaved data from raw audio. However, they do not quantify the importance of interleaving granularity for effective training. To study this, we construct two interleaving variants (see Fig. 2A)— (1) coarse interleaving, where we merge multiple consecutive diarized outputs into one if tagged with same speaker-ID, yielding long chunks, and (2) fine interleaving, where we keep all diarized outputs as is without merging, yielding short chunks. As expected, from Fig. 3, we find coarse interleaving leads to longer chunks (meanlength=19.2s) compared to fine interleaving (mean-length=5.2s). From Tab. 1, we note fine interleaving improves SQA performance by 3.1% on average, while matching text-only performance. This is a significant finding since the default approach in prior works (Ding et al., 2025; Li et al., 2025b) has been to merge same-speaker diarization outputs, yet our results advocate for more granular interleaving. Hence, for all our subsequent

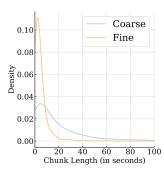


Figure 3: Audio chunk length distribution (in seconds) for our interleaving strategies.

experiments, we adopt fine interleaving for web-crawled speech-text pretraining by default.

Table 1: Fine chunking improves over coarse chunking.

Interleaving Granularity	Text Under	SQA (S→T) acc (%)				
	CoreEN	MMLU	SWQ	STQ	SLQ	Avg
Coarse	60.4	63.9	42.5	26.6	43.6	37.6
Fine	60.4	64.1	42.7	32.2	47.3	40.7

Takeaway: Fine-grained interleaving of speech-text pretraining data boosts SQA performance.

3.4 Constructing effective synthetic datasets

While web-crawled datasets offer massive volume, they often have poor domain coverage—their data distribution does not reflect the highest-priority domains for downstream deployment (Baack, 2024; Longpre et al., 2024). Often, sufficient data from many core domains simply does not exist or is hard to crawl (Zhang et al., 2024c; Fang et al., 2023b). Together, these motivate using synthetic data to augment existing data from web-crawls. Moreover, in our web-crawled audio data, we find noisy text-annotations (due to hallucinations from transcription models) and artifacts like background noise and speaker overlap. Thereby, we explore synthesizing interleaved speech-text datasets from existing text-only corpora. We build two synthetic datasets (see Fig. 2B) to augment web-crawled data-Knowledge-Rich Interleaved Speech-Text (Krist) and Question-Answering Speech-Text (Quest).

Knowledge-Rich Interleaved Speech-Text (Krist). We start from lightly-filtered web-crawled documents (similar to WARC files from CommonCrawl (2007)). We then apply URL-filtering to preserve knowledge-rich domains (list is in Appx. B.1). This is motivated by recent efforts advocating high-quality educational data for accelerating training (Penedo et al., 2024; Abdin et al., 2024; Gunasekar et al., 2023). Next, we use qpt-40-mini to extract and lightly rewrite the text-content from raw HTML, following Maini et al. (2024) (prompt used in Appx. B.2). We then segment the texts based on sentence-level splitting, to produce different text chunks. Finally, we synthesize audio for each chunk using melo-TTS (Zhao et al., 2023). To improve diversity, we randomly sample voices from 5 different accents. This pipeline yields \sim 4.6M hours of interleaved speech-text data.

Question-Answering Speech-Text (Quest). Since *Krist* is synthesized from HTML-extracted text, its constituent samples do not sound like natural conversations. We therefore build *Quest*, explicitly organized in a question-answering format to mimic real world audio. Starting from the same high-quality HTML pool as Krist, we first mine all possible question texts using regex-parsing. We then use gpt-4o to filter out invalid questions (examples in Appx. B.3). Finally, we use gpt-4o to generate responses along with a chain-of-thought (Wei et al., 2022) trace (prompts in Appx. B.2). We use the same sentence-level chunking as *Krist* to get \sim 0.9M hours of interleaved speech-text data.

Table 2: Synthetic speech-text interleaved data improves over web-crawl data.

Data Mix	Text Under	SQA (S \rightarrow T) acc (%)				
2 	CoreEN	MMLU	SWQ	STQ	SLQ	Avg
Web-crawl 100%	60.4	64.1	42.7	32.2	47.3	40.7
Web-crawl 53% + Krist 47%	60.8	64.8	43.4	29.2	52.0	41.5
Web-crawl 66% + Quest 34%	60.4	66.2	42.7	34.7	66.3	47.9
Web-crawl 59% + Quest 6% + Krist 35%	60.7	65.9	43.8	31.5	51.0	42.1
Web-crawl 40% + Quest 27% + Krist 33%	60.6	65.7	43.3	31.7	49.3	41.4

Results. We study the impact of mixing *Krist* and *Quest* with web-crawled data (mixed proportional to their approximate token counts) in Tab. 2. We find that mixing in *Krist* brings a 0.8% lift in SQA performance while also moderately benefitting text-only benchmarks, compared to training on web-crawl alone. Further, mixing *Quest* with web-crawl improves both MMLU and SQA performance by 2.1% and 7.2%. We hypothesize that question-answering format in interleaved training helps to sample-efficiently adapt to downstream SQA capabilities. We additionally explore two further ratios for mixing *Quest* and *Krist* with the web-crawled data—one where we sample according to approximate token-counts of each data source (59% web-crawl), and another where we upsample the synthetic proportion (40% web-crawl). Both settings improve over web-crawl by 1.4–0.7% SQA. However, due to complex interactions between mixing ratios and data repeats (Muennighoff et al., 2023), it is unclear how to construct an optimal mixture extracting the best of each data source (Shukor et al., 2025) (details on exact token counts in Appx. C). We leave this exploration for future work.

Takeaway: Synthetic datasets using TTS models bring gains when mixed with web-crawled data.

3.5 MODALITY SAMPLING SCHEMES FOR INTERLEAVED TRAINING

So far, we have discussed interleaved speech-text data *processing* and *curation* for improving SQA performance. However, we did not describe *how we sample modality chunks during interleaved training*. Here, we study two different sampling schemes as shown in Fig. 2C. Recollect that each interleaved training sample is of the form $X_i = \{(A_1, T_1), (A_2, T_2) \cdots (A_n, T_n)\}$.

Stochastic Sampling. In the first variant (used in prior experiments), at each chunk i, we randomly sample the chunk-modality with 0.5 probability. The modality sampling at each chunk i is independent of other chunks. We always start with an audio chunk A_1 , to ensure that there is at least 1 audio chunk.

Deterministic Sampling. While the stochastic variant allows flexibility potentially offering better generalization, it can restrict the number of *modality switches* during training. Hence, we test a deterministic approach, where we alternate between audio and text

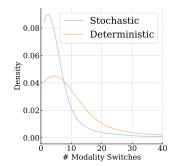


Figure 4: Number of modality switches during interleaved training for sampling schemes.

modalities at each chunk, i.e. we formulate the training sequence as $\{A_1, T_2, A_3 \cdots A_{n-1}, T_n\}$. This maximizes the number of modality switches for a given sample. Here too, we always start with A_1 .

Results. From Tab. 3, we find deterministic sampling boosts SQA performance by 1% on average over stochastic sampling. We posit that the number of modality switches during training affects the SQA performance—in Fig. 4, we plot the distribution of modality switches occurring during interleaved training, finding that stochastic sampling switches modalities quite infrequently, whereas

the deterministic approach has a higher number of switches during training. Indeed, the expected number of modality switches for a sample consisting of n chunks is n-1 for deterministic sampling and $\frac{n-1}{2}$ for stochastic sampling. By frequently switching modalities more often, deterministic sampling likely enables more effective cross-modal learning, thereby improving performance.

Table 3: Deterministic speech-text sampling improves over stochastic sampling.

Sampling scheme	Text Under	standing $(T \rightarrow T)$	SQA (S \rightarrow T) acc (%)				
Sumpring Series	CoreEN MMLU		SWQ	STQ	SLQ	Avg	
Stochastic Deterministic	60.6 60.1	65.7 65.2	43.3 44.2	31.7 31.2	49.3 51.7	41.4 42.4	

Takeaway: Deterministic sampling improves SQA ability over stochastic for interleaved training.

3.6 Our data-centric lessons transfer to understanding-only SpeechLMs

So far, we showed our three data-centric methods boost SQA significantly. These results were achieved while computing the loss on both audio and text tokens during interleaved training to support a native end-to-end speechLM. However, there is also great interest in developing an understanding-only speechLM, e.g. the Thinker model in Thinker-Talker architecture (Xu et al., 2025). In this vein, many prior works (Liu et al., 2025; Chu et al., 2024; Li et al., 2025b) apply loss masking on the audio tokens while doing speech-text interleaved training.

We hence test if our three data strategies also transfer to the audio-loss-masked setting. From Tab. 4, we find this indeed to be the case (9.3% average SQA lift). Further, we find absolute SQA performance improves significantly with loss-masking (51.8% with loss masking vs. 42.4% without). This result corroborates prior results (Liu et al., 2025; Li et al., 2025b; Chu et al., 2024) suggesting that, for small scale models there is an inherent modality conflict between audio and text tokens, which can lead to regressions when computing loss on both modalities.

Table 4: Our data-centric methods also work for understanding-only SpeechLM

Method	Text Under	rstanding $(T \rightarrow T)$	$SQA (S \rightarrow T) acc (\%)$			
	CoreEN	MMLU	SWQ	STQ	SLQ	Avg
Baseline (w/o loss-masking)	60.4	63.9	42.5	26.6	43.6	40.7
+ all data interventions	60.1	65.2	44.2	31.2	51.7	42.4
Baseline (w/ loss-masking)	61.7	66.5	45.9	34.0	47.7	42.5
+ all data interventions	61.8	67.3	45.7	44.6	65.0	51.8

Takeaway: Our three data interventions also transfer to the understanding-only SpeechLM setting.

4 Understanding Why Our Data Interventions Help

4.1 IMPROVED ALIGNMENT BETWEEN MODALITY DISTRIBUTIONS

Here, we aim to better understand why our data interventions (fine chunking + synthetic data mixing) improve over a baseline with coarse chunking and no synthetic data. One plausible hypothesis is that *fine interleaving and synthetic data close the gap between the model's audio-conditioned and text-conditioned output distribution*. Since we initialize from a well-trained language model, ensuring the audio-conditioned output distribution matches the distribution of text-conditioned outputs enables *strong modality alignment*. We now test if our data-centric approaches close this distribution gap.

Setup. We start with Spoken-LLaMA-Questions. For each test sample, we independently compute the token-wise teacher-forced probability distributions based on conditioning on audio and text questions separately. We then compute the mean token-wise KL-divergence values between the two probability distributions. For more details, definitions and other metrics, please refer to Appx. H.

Results. In Fig. 5, we plot the distribution of mean KL-divergence values between text-conditioned and audio-conditioned output distributions on the full Spoken-LLaMA-Questions test set. We find that fine interleaving induces lower KL-divergence values (mean=2.21) compared to coarse interleaving (mean=3.20). Moreover, a model trained with both fine interleaving and synthetic data further closes the distribution gap (mean KLD=1.47). This holds across other metrics and test sets too (see Appx. H). This result suggests our data interventions indeed help to close the gap between text- and audio-conditioned probability distributions, thereby better aligning the two modalities, leading to stronger SQA capability.

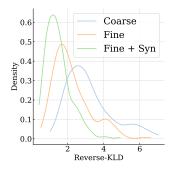


Figure 5: Our methods reduce the distribution gap between text and audio modalities.

4.2 SYNTHETIC DATA IMPROVES DOMAIN COVERAGE

In Sec. 3.4, we observed that our synthetic datasets improve both text and SQA performance significantly. Our central hypothesis for *why* is—*web-crawled data has a very skewed topic distribution and our synthetic data improves the domain coverage*. To help understand the composition of our web-crawled and synthetic datasets from a *topic* perspective, we leveraged the topic-domain classifier from (Wettig et al., 2025)¹, which can categorize texts in 24 different topic domains (an analysis with more fine-grained classifiers is in Appx. I.3). We run the classifier on 5000 random samples from each of our training datasets (*Web-crawl, Krist* and *Quest*). We also annotate topics covered in evaluation datasets. From our results in Fig. 6 (more results in Appx. I), we make two key observations:

- Web-crawled data is highly skewed and is majorly comprised of *entertainment*, *sports* and *fitness*, *religion* and *social life* domains. This is not surprising given that most of our web-crawled data is sourced from podcasts, interviews, talk-shows and monologues.
- **Synthetic data improves topic coverage.** It is evident that both *Krist* and *Quest* oversample data from the domains of *science and tech*, *health*, *education and jobs*, and *finance*, all of which are extremely under-represented in the web-crawled data.

Therefore, by enabling broader coverage of topic domains, our synthetic datasets help to (i) close the distribution mismatch between the raw web-crawled data and the downstream evaluation datasets, and (ii) enhance the diversity of our pretraining data distribution. Prior work in the language space has also discussed the importance of training data diversity and domain coverage (Wettig et al., 2025; Nguyen et al., 2025a; Maini et al., 2025; Wang et al., 2025b; Maini et al., 2024).

4.3 Analysing train-test contamination

Given the significant boosts induced by our synthetic datasets, a natural question arises—*Is there test-set leakage, and if so, how does it impact SQA performance?* To address this, we conduct a contamination analysis with two goals in mind: (1) identify the proportion of test samples that are likely contaminated in our training data, and (2) understand the performance impact of this leakage.

Contamination detection. To find extent of contamination in our synthetic datasets, we follow recent works (Singh et al., 2024; Sainz et al., 2024; Li et al., 2024; Dubey et al., 2024) and use n-gram token overlaps. While prior works used n=13, we opt for a window from n=6 to n=13 to improve recall, at the expense of more false-positives. We use the gpt-40 tokenizer and apply lower-case normalization pre-tokenizing. We mark a test sample as contaminated if we find a matching n-gram in any equivalent n-token span of a synthetic dataset. We consider all three SQA test sets for analysis, and concatenate the question and answer of each sample for matching. For train sets, we take samples from the original seed text-datasets (from which we synthesize audio) for detecting matches.

Proportion of contamination. We report the proportion of contaminated test samples in Tab. 11. We find that the *Quest* dataset has almost no contamination, while *Krist* has a small, yet nonnegligible amount of contamination. Overall, SWQ eval is barely contaminated (0.4%) while STQ and SLQ evals have 2.5% and 7.7% contaminated samples respectively. Importantly, note that due to our windowed n-gram approach, we have many false-positive matches (examples of matches are in Appx. J.1). However, we keep all matches to be as conservative as possible while analysing

¹https://huggingface.co/WebOrganizer/TopicClassifier-NoURL

Distributions of topic domains in our evaluation and training datasets

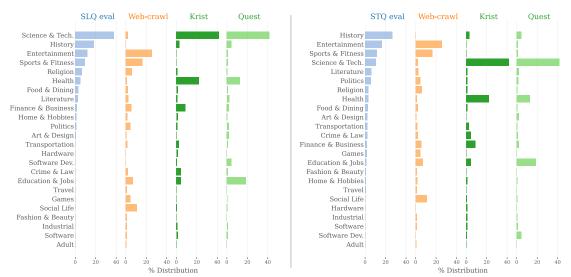


Figure 6: **Synthetic data improves domain coverage.** We plot the distribution of topic domains in our evaluation datasets (in blue, Spoken-LLaMA-Questions (*left*) and Spoken-TriviaQA (*right*)) and contrast them with the topic distribution of web-crawled and synthetic datasets. Our synthetic datasets (*Krist* and *Quest*) fill gaps in domains that are under-represented in the web-crawled data, reducing distribution mismatch, thereby enabling better downstream SQA performance.

effect of contamination. To understand the impact of test-set contamination on downstream model performance, we consider the tests sets with these contaminated samples removed as *clean* sets—*SWQ-clean* has 996 samples, *STQ-clean* has 975 samples, and *SLQ-clean* has 277 samples.

Significance testing setup. We conduct a one-sided significance test on the differences between performance on the full test set (including all contaminated samples) and performance on the clean set (removing all contaminated samples). To control for the accuracy difference induced by reducing test set size for the clean sets, we compute the random removal baseline accuracy—model performance after removing the same number of randomly selected test samples, averaged across 100 bootstrap replicates with different seeds. We compute empirical p-values by comparing clean accuracy against the bootstrap distribution. Under this setting, our null hypothesis is: observed model accuracy on the full test set is not artificially inflated by contamination. For more details on the significance testing setup, refer Appx. J.3.

Results. We apply the previously described significance testing procedure for all 4 models trained in Sec. 3.4 that use synthetic data in their data mixture. In Fig. 7, we plot the absolute differences between the *clean* and *random removal mean* accuracy for all eval-model pairs.

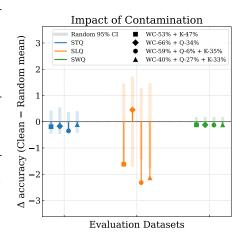


Figure 7: We plot differences between clean and random removal mean accuracies along with 95% CIs—our results suggest contamination has minor effect.

We find that contamination does not seem to significantly improve performance for *Spoken-TriviaQA* and *Spoken-Web-Questions* evaluations. For *Spoken-LLaMA-Questions*, contamination seems to be having a minor effect (1.4-2.1%) when *Krist* is in the data mixture. However, we find that the effect of contamination is not statistically significant in almost all the settings (at a significance level of α =0.01). We provide a more in-depth analysis with confidence intervals (CIs) and the *p*-values in Appx. J.3. Additionally, we note that the performance boosts on Spoken-LLaMA-Questions due to synthetic data observed in Sec. 3.4 (3.7%-19%) far exceed the clean vs random-removal-mean

accuracy differences observed here (upto 2%). Taken together, these results suggest that test-set contamination does not play a major role in explaining the accuracy boosts observed in Sec. 3.4.

5 Spelangy: Bringing it all together

Equipped with our key data-centric insights from previous sections, we now train a 3.8B SpeechLM, called **SpeLangy**. We use the same configuration as before, with 16, 384 sequence length trained for 1.67T speech-text tokens. We compare against Kimi-Audio (Ding et al., 2025), Qwen-Audio (Chu et al., 2023), and Qwen2-Audio (Chu et al., 2024). We additionally compare two post-trained models—Voxtral-mini (Liu et al., 2025) and GLM-4-Voice (Zeng et al., 2024a)—with the caveat that having undergone instruction-tuning, they are not directly comparable to base models (Dominguez-Olmedo et al., 2024). To ensure our training recipe does not degrade language performance, we also compare against strong open-weights base language models on standard text-only benchmarks.

Results. From Tab. 5 we find **SpeLangy** outperforms Kimi-Audio, Qwen-Audio and Qwen-2-Audio by 10.2%, 11.1% and 9.8% on average across the three SQA benchmarks, while being $2.8 \times$, $2.2 \times$ and $2.2 \times$ smaller in size. Further, we obtain competitive performance with the strongly post-trained Voxtral-mini and GLM-4-Voice, without having undergone any specific instruction-tuning. In Tab. 6, we compare text performance of **SpeLangy** with the base LM that we initialize from—we observe large boosts across the board compared to the base-LM. Further, our model is competitive with Gemma-2 (Team et al., 2024), Gemma-3 (Team et al., 2025) and Qwen-2.5 (Yang et al., 2024) models, all of which are leading open-weights models, highlighting the strength of our **SpeLangy**.

Table 5: **Spoken Question-Answering** ($S \rightarrow T$) **comparison.** We report results for SoTA SpeechLMs and **SpeLangy**. Where possible, we report results obtained using the pretrained base models (if no base models are publicly released, we evaluate the post-trained checkpoints and make a note of this).

Type	Model	# Params	SWQ	STQ	SLQ	Average
Base	Kimi-Audio Qwen-Audio Qwen-2-Audio	10.5B 8.4B 8.4B	44.0 45.7 45.7	33.8 30.3 33.4	47.0 46.0 47.0	41.6 40.7 42.0
	SpeLangy	3.8B	45.7	44.6	65.0	51.8
SFT	Voxtral-mini GLM-4-Voice	4.7B 9.9B	41.6 43.3	46.6 52.4	65.3 64.7	51.2 53.4

Table 6: **Text Understanding** ($T \rightarrow T$) **comparison.** We compare **SpeLangy** with leading text-only base models of same size-class. *Text-init* is the model we start continued-pretraining from. Our model is competitive with all compared models, highlighting that we strongly preserve text-only capabilities.

Model	# Params	CoreEN	MMLU	GSM8k	HumanEval
Text-init	3B	62.4	62.2	47.1	29.9
Gemma-2 Gemma-3 Qwen-2.5	2.6B 4B 3B	- - -	56.1 62.8 <u>65.6</u>	30.3 38.4 79.1	19.5 36.0 42.1
SpeLangy	3.8B	61.8	67.3	<u>71.9</u>	<u>37.6</u>

6 Conclusion

In this work, we studied three data-curation methods for speech-language interleaved pretraining to enhance spoken question-answering (SQA) capabilities. We found fine-grained interleaving of speech-text chunks bringing large gains, while synthetic datasets synthesized from knowledge-rich seed text-datasets also boosted performance. Deterministic sampling of speech-text chunks during interleaved pretraining further improved SQA results. We showed that these data-centric recipes strengthen alignment between the speech and text modalities and broaden domain coverage of pretraining datasets. Distilling these insights, we pretrained **SpeLangy**, achieving competitive performance with larger models. We hope our insights motivate more data-centric SpeechLM work.

REFERENCES

- Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*, 2024.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Robin Algayres, Yossi Adi, Tu Anh Nguyen, Jade Copet, Gabriel Synnaeve, Benoit Sagot, and Emmanuel Dupoux. Generative spoken language model based on continuous word-sized audio tokens. *arXiv preprint arXiv:2310.05224*, 2023.
- Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav, et al. Smollm2: When smol goes big–data-centric training of a small language model. *arXiv* preprint arXiv:2502.02737, 2025.
- Siddhant Arora, Kai-Wei Chang, Chung-Ming Chien, Yifan Peng, Haibin Wu, Yossi Adi, Emmanuel Dupoux, Hung-Yi Lee, Karen Livescu, and Shinji Watanabe. On the landscape of spoken language models: A comprehensive survey. *arXiv* preprint arXiv:2504.08528, 2025.
- Stefan Baack. A critical analysis of the largest source for generative ai training data: Common crawl. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pp. 2199–2208, 2024.
- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on freebase from question-answer pairs. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp. 1533–1544, 2013.
- Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.
- Sebastian Bordt, Suraj Srinivas, Valentyn Boreiko, and Ulrike von Luxburg. How much can we forget about data contamination? *arXiv preprint arXiv:2410.03249*, 2024.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, et al. Jax: Autograd and xla. *Astrophysics Source Code Library*, pp. ascl–2111, 2021.
- Hervé Bredin. pyannote. audio 2.1 speaker diarization pipeline: principle, benchmark, and recipe. In 24th INTERSPEECH Conference (INTERSPEECH 2023), pp. 1983–1987. ISCA, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021.

- Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv* preprint arXiv:2311.07919, 2023.
 - Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, et al. Qwen2-audio technical report. *arXiv preprint arXiv:2407.10759*, 2024.
 - Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
 - Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
 - Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. arXiv preprint arXiv:2507.06261, 2025.
 - CommonCrawl. Common Crawl: Open Repository of Web Crawl Data. https://registry.opendata.aws/commoncrawl/, 2007. Accessed: 2025-09-15.
 - Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard Grave, and Neil Zeghidour. Moshi: a speech-text foundation model for real-time dialogue. *arXiv preprint arXiv:2410.00037*, 2024.
 - Ding Ding, Zeqian Ju, Yichong Leng, Songxiang Liu, Tong Liu, Zeyu Shang, Kai Shen, Wei Song, Xu Tan, Heyi Tang, et al. Kimi-audio technical report. *arXiv preprint arXiv:2504.18425*, 2025.
 - Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. Documenting large webtext corpora: A case study on the colossal clean crawled corpus. *arXiv preprint arXiv:2104.08758*, 2021.
 - Ricardo Dominguez-Olmedo, Florian E Dorner, and Moritz Hardt. Training on the test task confounds evaluation and emergence. *arXiv preprint arXiv:2407.07890*, 2024.
 - Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.
 - Yanai Elazar, Akshita Bhagia, Ian Magnusson, Abhilasha Ravichander, Dustin Schwenk, Alane Suhr, Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, et al. What's in my big data? *arXiv* preprint arXiv:2310.20707, 2023.
 - David Fan, Shengbang Tong, Jiachen Zhu, Koustuv Sinha, Zhuang Liu, Xinlei Chen, Michael Rabbat, Nicolas Ballas, Yann LeCun, Amir Bar, et al. Scaling language-free visual representation learning. *arXiv preprint arXiv:2504.01017*, 2025.
 - Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander Toshev, and Vaishaal Shankar. Data filtering networks. *arXiv preprint arXiv:2309.17425*, 2023a.
 - Alex Fang, Simon Kornblith, and Ludwig Schmidt. Does progress on imagenet transfer to real-world datasets? *Advances in Neural Information Processing Systems*, 36:25050–25080, 2023b.
- Enrico Fini, Mustafa Shukor, Xiujun Li, Philipp Dufter, Michal Klein, David Haldimann, Sai Aitharaju, Victor G Turrisi da Costa, Louis Béthune, Zhe Gan, et al. Multimodal autoregressive pretraining of large vision encoders. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 9641–9654, 2025.
 - Jonathan G Fiscus. A post-processing system to yield reduced word error rates: Recognizer output voting error reduction (rover). In 1997 IEEE Workshop on Automatic Speech Recognition and Understanding Proceedings, pp. 347–354. IEEE, 1997.

- Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen,
 Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, et al. Datacomp: In search of the
 next generation of multimodal datasets. *Advances in Neural Information Processing Systems*, 36:
 27092–27112, 2023.
 - Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv* preprint arXiv:2101.00027, 2020.
 - Xuelong Geng, Kun Wei, Qijie Shao, Shuiyun Liu, Zhennan Lin, Zhixian Zhao, Guojian Li, Wenjie Tian, Peikun Chen, Yangze Li, et al. Osum: Advancing open speech understanding models with limited resources in academia. *arXiv preprint arXiv:2501.13306*, 2025.
 - Adhiraj Ghosh, Sebastian Dziadzio, Ameya Prabhu, Vishaal Udandarao, Samuel Albanie, and Matthias Bethge. Onebench to test them all: Sample-level benchmarking over open-ended capabilities. *arXiv preprint arXiv:2412.06745*, 2024.
 - Sreyan Ghosh, Zhifeng Kong, Sonal Kumar, S Sakshi, Jaehyeon Kim, Wei Ping, Rafael Valle, Dinesh Manocha, and Bryan Catanzaro. Audio flamingo 2: An audio-language model with long-audio understanding and expert reasoning abilities. *arXiv preprint arXiv:2503.03983*, 2025.
 - Arushi Goel, Sreyan Ghosh, Jaehyeon Kim, Sonal Kumar, Zhifeng Kong, Sang-gil Lee, Chao-Han Huck Yang, Ramani Duraiswami, Dinesh Manocha, Rafael Valle, et al. Audio flamingo 3: Advancing audio intelligence with fully open large audio language models. *arXiv preprint arXiv:2507.08128*, 2025.
 - Yuan Gong, Alexander H Liu, Hongyin Luo, Leonid Karlinsky, and James Glass. Joint audio and speech understanding. In 2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU), pp. 1–8. IEEE, 2023.
 - Yuling Gu, Oyvind Tafjord, Bailey Kuehl, Dany Haddad, Jesse Dodge, and Hannaneh Hajishirzi. Olmes: A standard for language model evaluations. *arXiv preprint arXiv:2406.08446*, 2024.
 - Etash Guha, Ryan Marten, Sedrick Keh, Negin Raoof, Georgios Smyrnis, Hritik Bansal, Marianna Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, et al. Openthoughts: Data recipes for reasoning models. *arXiv preprint arXiv:2506.04178*, 2025.
 - Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. Conformer: Convolution-augmented transformer for speech recognition. *arXiv* preprint arXiv:2005.08100, 2020.
 - Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*, 2023.
 - Tom Gunter, Zirui Wang, Chong Wang, Ruoming Pang, Andy Narayanan, Aonan Zhang, Bowen Zhang, Chen Chen, Chung-Cheng Chiu, David Qiu, et al. Apple intelligence foundation language models. *arXiv preprint arXiv:2407.21075*, 2024.
 - Michael Hassid, Tal Remez, Tu Anh Nguyen, Itai Gat, Alexis Conneau, Felix Kreuk, Jade Copet, Alexandre Defossez, Gabriel Synnaeve, Emmanuel Dupoux, et al. Textually pretrained speech language models. *Advances in Neural Information Processing Systems*, 36:63483–63501, 2023.
 - Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
 - Andreas Hochlehnert, Hardik Bhatnagar, Vishaal Udandarao, Samuel Albanie, Ameya Prabhu, and Matthias Bethge. A sober look at progress in language model reasoning: Pitfalls and paths to reproducibility. *arXiv preprint arXiv:2504.07086*, 2025.
 - Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.

- Shahab Jalalvand, Matteo Negri, Daniele Falavigna, and Marco Turchi. Driving rover with segment-based asr quality estimation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 1095–1105, 2015.
 - Minhao Jiang, Ken Ziyu Liu, Ming Zhong, Rylan Schaeffer, Siru Ouyang, Jiawei Han, and Sanmi Koyejo. Investigating data contamination for pre-training language models. *arXiv* preprint *arXiv*:2401.06059, 2024.
 - Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.
 - Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. Large language models struggle to learn long-tail knowledge. In *International conference on machine learning*, pp. 15696–15707. PMLR, 2023.
 - Nikhil Kandpal, Brian Lester, Colin Raffel, Sebastian Majstorovic, Stella Biderman, Baber Abbasi, Luca Soldaini, Enrico Shippole, A Feder Cooper, Aviya Skowron, et al. The common pile v0. 1: An 8tb dataset of public domain and openly licensed text. *arXiv preprint arXiv:2506.05209*, 2025.
 - Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv* preprint arXiv:2001.08361, 2020.
 - Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping, Rafael Valle, and Bryan Catanzaro. Audio flamingo: A novel audio language model with few-shot learning and dialogue abilities. *arXiv* preprint arXiv:2402.01831, 2024.
 - Taku Kudo and John Richardson. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv* preprint arXiv:1808.06226, 2018.
 - Kushal Lakhotia, Eugene Kharitonov, Wei-Ning Hsu, Yossi Adi, Adam Polyak, Benjamin Bolte, Tu-Anh Nguyen, Jade Copet, Alexei Baevski, Abdelrahman Mohamed, et al. On generative spoken language modeling from raw audio. *Transactions of the Association for Computational Linguistics*, 9:1336–1354, 2021.
 - Mark Lee, Tom Gunter, Chang Lan, John Peebles, Hanzhi Zhou, Kelvin Zou, Sneha Bangalore, Chung-Cheng Chiu, Nan Du, Xianzhi Du, et al. Axlearn: Modular large model training on heterogeneous infrastructure. *arXiv preprint arXiv:2507.05411*, 2025.
 - Ethan Li, Anders Boesen Lindbo Larsen, Chen Zhang, Xiyou Zhou, Jun Qin, Dian Ang Yap, Narendran Raghavan, Xuankai Chang, Margit Bowler, Eray Yildiz, et al. Apple intelligence foundation language models: Tech report 2025. *arXiv preprint arXiv:2507.13575*, 2025a.
 - Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik Bansal, Etash Guha, Sedrick Scott Keh, Kushal Arora, et al. Datacomp-lm: In search of the next generation of training sets for language models. *Advances in Neural Information Processing Systems*, 37:14200–14282, 2024.
 - Tianpeng Li, Jun Liu, Tao Zhang, Yuanbo Fang, Da Pan, Mingrui Wang, Zheng Liang, Zehuan Li, Mingan Lin, Guosheng Dong, et al. Baichuan-audio: A unified framework for end-to-end speech interaction. *arXiv preprint arXiv:2502.17239*, 2025b.
 - Yang Li, Youssef Emad, Karthik Padthe, Jack Lanchantin, Weizhe Yuan, Thao Nguyen, Jason Weston, Shang-Wen Li, Dong Wang, Ilia Kulikov, et al. Naturalthoughts: Selecting and distilling reasoning traces for general reasoning tasks. *arXiv preprint arXiv:2507.01921*, 2025c.
 - Alexander H Liu, Andy Ehrenberg, Andy Lo, Clément Denoix, Corentin Barreau, Guillaume Lample, Jean-Malo Delignon, Khyathi Raghavi Chandu, Patrick von Platen, Pavankumar Reddy Muddireddy, et al. Voxtral. *arXiv preprint arXiv:2507.13264*, 2025.

- Shayne Longpre, Robert Mahari, Ariel Lee, Campbell Lund, Hamidah Oderinwale, William Brannon, Nayan Saxena, Naana Obeng-Marnu, Tobin South, Cole Hunter, et al. Consent in crisis: The rapid decline of the ai data commons. *Advances in Neural Information Processing Systems*, 37: 108042–108087, 2024.
 - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint* arXiv:1711.05101, 2017.
 - Pratyush Maini, Skyler Seto, He Bai, David Grangier, Yizhe Zhang, and Navdeep Jaitly. Rephrasing the web: A recipe for compute and data-efficient language modeling. *arXiv* preprint *arXiv*:2401.16380, 2024.
 - Pratyush Maini, Vineeth Dorna, Parth Doshi, Aldo Carranza, Fan Pan, Jack Urbanek, Paul Burstein, Alex Fang, Alvin Deng, Amro Abbas, et al. Beyondweb: Lessons from scaling synthetic data for trillion-scale pretraining. *arXiv* preprint arXiv:2508.10975, 2025.
 - Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruti Shah, Xianzhi Du, Futang Peng, Anton Belyi, et al. Mm1: methods, analysis and insights from multimodal llm pre-training. In *European Conference on Computer Vision*, pp. 304–323. Springer, 2024.
 - Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschannen. Finite scalar quantization: Vq-vae made simple. *arXiv preprint arXiv:2309.15505*, 2023.
 - David Mizrahi, Anders Boesen Lindbo Larsen, Jesse Allardice, Suzie Petryk, Yuri Gorokhov, Jeffrey Li, Alex Fang, Josh Gardner, Tom Gunter, and Afshin Dehghan. Language models improve when pretraining data matches target tasks. *arXiv preprint arXiv:2507.12466*, 2025.
 - Hossein Mobahi, Mehrdad Farajtabar, and Peter Bartlett. Self-distillation amplifies regularization in hilbert space. *Advances in Neural Information Processing Systems*, 33:3351–3361, 2020.
 - Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language models. *Advances in Neural Information Processing Systems*, 36:50358–50376, 2023.
 - Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.
 - Eliya Nachmani, Alon Levkovitch, Roy Hirsch, Julian Salazar, Chulayuth Asawaroengchai, Soroosh Mariooryad, Ehud Rivlin, RJ Skerry-Ryan, and Michelle Tadmor Ramanovich. Spoken question answering and speech continuation using spectrogram-powered llm. *arXiv preprint arXiv:2305.15255*, 2023.
 - Huong Ngo, Matt Deitke, Martijn Bartelds, Sarah Pratt, Josh Gardner, Matt Jordan, and Ludwig Schmidt. Olmoasr: Open models and data for training robust speech recognition models. *arXiv* preprint arXiv:2508.20869, 2025.
 - Thao Nguyen, Yang Li, Olga Golovneva, Luke Zettlemoyer, Sewoong Oh, Ludwig Schmidt, and Xian Li. Recycling the web: A method to enhance pre-training data quality and quantity for language models. *arXiv preprint arXiv:2506.04689*, 2025a.
 - Tu Anh Nguyen, Benjamin Muller, Bokai Yu, Marta R Costa-Jussa, Maha Elbayad, Sravya Popuri, Christophe Ropers, Paul-Ambroise Duquenne, Robin Algayres, Ruslan Mavlyutov, et al. Spirit-lm: Interleaved spoken and written language model. *Transactions of the Association for Computational Linguistics*, 13:30–52, 2025b.
 - Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*, 2024.
 - OpenAI. Openai o3 and o4-mini system card, 2024. URL https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf. Accessed: 2025-05-11.

- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
 - Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset: Word prediction requiring a broad discourse context. *arXiv preprint arXiv:1606.06031*, 2016.
 - Shubham Parashar, Zhiqiu Lin, Tian Liu, Xiangjue Dong, Yanan Li, Deva Ramanan, James Caverlee, and Shu Kong. The neglected tails in vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12988–12997, 2024.
 - Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv* preprint arXiv:2306.01116, 2023.
 - Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin A Raffel, Leandro Von Werra, Thomas Wolf, et al. The fineweb datasets: Decanting the web for the finest text data at scale. *Advances in Neural Information Processing Systems*, 37:30811–30849, 2024.
 - Yifan Peng, Shakeel Muhammad, Yui Sudo, William Chen, Jinchuan Tian, Chyi-Jiunn Lin, and Shinji Watanabe. Owsm v4: Improving open whisper-style speech models via data scaling and cleaning. arXiv preprint arXiv:2506.00338, 2025.
 - Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
 - Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pp. 28492–28518. PMLR, 2023.
 - Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
 - Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
 - Ankit Singh Rawat, Veeranjaneyulu Sadhanala, Afshin Rostamizadeh, Ayan Chakrabarti, Wittawat Jitkrittum, Vladimir Feinberg, Seungyeon Kim, Hrayr Harutyunyan, Nikunj Saunshi, Zachary Nado, et al. A little help goes a long way: Efficient llm training by leveraging small lms. *arXiv* preprint arXiv:2410.18779, 2024.
 - Marvin Ritter, Ihor Indyk, Aayush Singh, Andrew Audibert, Anoosha Seelam, Camelia Hanes, Eric Lau, Jacek Olesiak, Jiyang Kang, and Xihui Wu. Grain feeding jax models, 2023. URL http://github.com/google/grain.
 - Paul K Rubenstein, Chulayuth Asawaroengchai, Duc Dung Nguyen, Ankur Bapna, Zalán Borsos, Félix de Chaumont Quitry, Peter Chen, Dalia El Badawy, Wei Han, Eugene Kharitonov, et al. Audiopalm: A large language model that can speak and listen. *arXiv preprint arXiv:2306.12925*, 2023.
 - Noveen Sachdeva and Julian McAuley. Data distillation: A survey. *arXiv preprint arXiv:2301.04272*, 2023.
 - Oscar Sainz, Iker García-Ferrero, Alon Jacovi, Jon Ander Campos, Yanai Elazar, Eneko Agirre, Yoav Goldberg, Wei-Lin Chen, Jenny Chim, Leshem Choshen, et al. Data contamination report from the 2024 conda shared task. *arXiv preprint arXiv:2407.21530*, 2024.
 - Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.

- Dvir Samuel, Rami Ben-Ari, Simon Raviv, Nir Darshan, and Gal Chechik. Generating images of rare concepts using pre-trained diffusion models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 4695–4703, 2024.
 - Mustafa Shukor, Enrico Fini, Victor Guilherme Turrisi da Costa, Matthieu Cord, Joshua Susskind, and Alaaeldin El-Nouby. Scaling laws for native multimodal models. *arXiv preprint arXiv:2504.07951*, 2025.
 - Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Dinov3. *arXiv* preprint arXiv:2508.10104, 2025.
 - Aaditya K Singh, Muhammed Yusuf Kocyigit, Andrew Poulton, David Esiobu, Maria Lomeli, Gergely Szilvasy, and Dieuwke Hupkes. Evaluation data contamination in llms: how do we measure it and (when) does it matter? *arXiv preprint arXiv:2411.03923*, 2024.
 - Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. Dolma: An open corpus of three trillion tokens for language model pretraining research. *arXiv preprint arXiv:2402.00159*, 2024.
 - Saurabh Srivastava, Anto PV, Shashank Menon, Ajay Sukumar, Alan Philipose, Stevin Prince, Sooraj Thomas, et al. Functional benchmarks for robust evaluation of reasoning performance, and the reasoning gap. *arXiv preprint arXiv:2402.19450*, 2024.
 - Dan Su, Kezhi Kong, Ying Lin, Joseph Jennings, Brandon Norick, Markus Kliegl, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. Nemotron-cc: Transforming common crawl into a refined long-horizon pretraining dataset. *arXiv* preprint arXiv:2412.02595, 2024.
 - Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and Chao Zhang. Salmonn: Towards generic hearing abilities for large language models. *arXiv preprint arXiv:2310.13289*, 2023.
 - Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
 - Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025.
 - Jinchuan Tian, Yifan Peng, William Chen, Kwanghee Choi, Karen Livescu, and Shinji Watanabe. On the effects of heterogeneous data sources on speech-to-text foundation models. *arXiv preprint arXiv:2406.09282*, 2024.
 - Peter Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Adithya Jairam Vedagiri IYER, Sai Charitha Akula, Shusheng Yang, Jihan Yang, Manoj Middepogu, Ziteng Wang, et al. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *Advances in Neural Information Processing Systems*, 37:87310–87356, 2024.
 - Trieu H Trinh and Quoc V Le. A simple method for commonsense reasoning. *arXiv preprint* arXiv:1806.02847, 2018.
 - Yuan Tseng, Titouan Parcollet, Rogier van Dalen, Shucong Zhang, and Sourav Bhattacharya. Evaluation of llms in speech is often flawed: Test set contamination in large language models for speech recognition. *arXiv preprint arXiv:2505.22251*, 2025.
 - Vishaal Udandarao, Ameya Prabhu, Adhiraj Ghosh, Yash Sharma, Philip Torr, Adel Bibi, Samuel Albanie, and Matthias Bethge. No" zero-shot" without exponential data: Pretraining concept frequency determines multimodal model performance. *Advances in Neural Information Processing Systems*, 37:61735–61792, 2024.

- Vishaal Udandarao, Nikhil Parthasarathy, Muhammad Ferjad Naeem, Talfan Evans, Samuel Albanie, Federico Tombari, Yongqin Xian, Alessio Tonioni, and Olivier J Hénaff. Active data curation effectively distills large-scale multimodal models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 14422–14437, 2025.
- Congchao Wang, Sean Augenstein, Keith Rush, Wittawat Jitkrittum, Harikrishna Narasimhan, Ankit Singh Rawat, Aditya Krishna Menon, and Alec Go. Cascade-aware training of language models. *arXiv preprint arXiv:2406.00060*, 2024.
- Enzhi Wang, Qicheng Li, Zhiyuan Tang, and Yuhang Jia. Cross-modal knowledge distillation for speech large language models. *arXiv preprint arXiv:2509.14930*, 2025a.
- Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation. *arXiv* preprint arXiv:1811.10959, 2018.
- Zaitian Wang, Jinghan Zhang, Xinhao Zhang, Kunpeng Liu, Pengfei Wang, and Yuanchun Zhou. Diversity-oriented data augmentation with large language models. *arXiv preprint arXiv:2502.11671*, 2025b.
- Maurice Weber, Dan Fu, Quentin Anthony, Yonatan Oren, Shane Adams, Anton Alexandrov, Xiaozhong Lyu, Huu Nguyen, Xiaozhe Yao, Virginia Adams, et al. Redpajama: an open dataset for training large language models. *Advances in neural information processing systems*, 37: 116462–116492, 2024.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Johannes Welbl, Nelson F Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. *arXiv* preprint arXiv:1707.06209, 2017.
- Alexander Wettig, Kyle Lo, Sewon Min, Hannaneh Hajishirzi, Danqi Chen, and Luca Soldaini. Organize the web: Constructing domains enhances pre-training data curation. *arXiv* preprint *arXiv*:2502.10341, 2025.
- Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Ben Feuer, Siddhartha Jain, Ravid Shwartz-Ziv, Neel Jain, Khalid Saifullah, Siddartha Naidu, et al. Livebench: A challenging, contamination-free llm benchmark. *arXiv preprint arXiv:2406.19314*, 4, 2024.
- Thaddäus Wiedemer, Yash Sharma, Ameya Prabhu, Matthias Bethge, and Wieland Brendel. Pretraining frequency predicts compositional generalization of clip on real-world tasks. *arXiv preprint arXiv:2502.18326*, 2025.
- Jack Wildman, Nikos I Bosse, Daniel Hnyk, Peter Mühlbacher, Finn Hambly, Jon Evans, Dan Schwarz, Lawrence Phillips, et al. Bench to the future: A pastcasting benchmark for forecasting agents. *arXiv preprint arXiv:2506.21558*, 2025.
- Ronald J Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280, 1989.
- Boyong Wu, Chao Yan, Chen Hu, Cheng Yi, Chengli Feng, Fei Tian, Feiyu Shen, Gang Yu, Haoyang Zhang, Jingbei Li, et al. Step-audio 2 technical report. *arXiv preprint arXiv:2507.16632*, 2025.
- LLM-Core-Team Xiaomi. Mimo-audio: Audio language models are few-shot learners, 2025. URL https://github.com/XiaomiMiMo/MiMo-Audio.
- Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, et al. Qwen2. 5-omni technical report. *arXiv preprint arXiv:2503.20215*, 2025.
- An Yang et al. Qwen2.5 technical report. arXiv preprint arXiv:2412.15115, 2024.
- Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E Gonzalez, and Ion Stoica. Rethinking benchmark and contamination for language models with rephrased samples. *arXiv* preprint arXiv:2311.04850, 2023.

- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
 - Aohan Zeng, Zhengxiao Du, Mingdao Liu, Kedong Wang, Shengmin Jiang, Lei Zhao, Yuxiao Dong, and Jie Tang. Glm-4-voice: Towards intelligent and human-like end-to-end spoken chatbot. *arXiv* preprint arXiv:2412.02612, 2024a.
 - Aohan Zeng, Zhengxiao Du, Mingdao Liu, Lei Zhang, Shengmin Jiang, Yuxiao Dong, and Jie Tang. Scaling speech-text pre-training with synthetic interleaved data. *arXiv preprint arXiv:2411.17607*, 2024b.
 - Zhiyuan Zeng, Jiashuo Liu, Siyuan Chen, Tianci He, Yali Liao, Jinpeng Wang, Zaiyuan Wang, Yang Yang, Lingyue Yin, Mingren Yin, et al. Futurex: An advanced live benchmark for llm agents in future prediction. *arXiv preprint arXiv:2508.11987*, 2025.
 - Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12104–12113, 2022.
 - Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities. *arXiv* preprint arXiv:2305.11000, 2023.
 - Ge Zhang et al. Finefineweb: A comprehensive study on fine-grained domain web corpus, December 2024a. URL [https://huggingface.co/datasets/m-a-p/FineFineWeb] (https://huggingface.co/datasets/m-a-p/FineFineWeb).
 - Jieyu Zhang, Weikai Huang, Zixian Ma, Oscar Michel, Dong He, Tanmay Gupta, Wei-Chiu Ma, Ali Farhadi, Aniruddha Kembhavi, and Ranjay Krishna. Task me anything. *arXiv preprint arXiv:2406.11775*, 2024b.
 - Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 3713–3722, 2019.
 - Linfeng Zhang, Chenglong Bao, and Kaisheng Ma. Self-distillation: Towards efficient and compact neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(8):4388–4403, 2021.
 - Yunhua Zhang, Hazel Doughty, and Cees GM Snoek. Low-resource vision challenges for foundation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21956–21966, 2024c.
 - Dora Zhao, Jerone TA Andrews, Orestis Papakyriakopoulos, and Alice Xiang. Position: measure dataset diversity, don't just claim it. *arXiv preprint arXiv:2407.08188*, 2024.
 - Wenliang Zhao, Xumin Yu, and Zengyi Qin. Melotts: High-quality multi-lingual multi-accent text-to-speech, 2023. URL https://github.com/myshell-ai/MeloTTS.

A PREPROCESSING WEB-CRAWLED AUDIO AS INTERLEAVED TRAINING DATA

In this section, we provide more details about each step in our data processing pipeline for converting web-crawled audio into interleaved speech-text format. We highlight all the components in Fig. 8.

Processing web-crawled audio into speech-text interleaved data

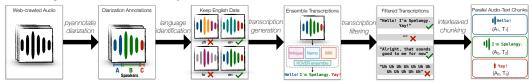


Figure 8: Our processing pipeline to convert raw web-crawled audio into trainable speech-text data.

Raw Audio. We start with a large corpus (>10M hours) of conversational-speech audios crawled from the web. These audios are sourced from a range of web domains, filtered to remove other audio types like music, ads and background audios. Our audio corpus primarily consist of podcasts, interviews and monologue speeches.

Speaker Diarization. Our first processing step involves identifying different speakers in each audio sample. We use pyannotate (Bredin, 2023) to annotate each audio sample into speaker diarized outputs. For each audio, the diarization procedure outputs a list of (audio-start, audio-end, speakerID) triplets. An example of a diarization output on an audio sample is shown below:

```
[{'start': 0.031, 'end': 5.971, 'speaker': 'SPEAKER_06'},
{'start': 7.085, 'end': 10.493, 'speaker': 'SPEAKER_06'},
{'start': 11.607, 'end': 13.278, 'speaker': 'SPEAKER_06'},
{'start': 13.565, 'end': 16.315, 'speaker': 'SPEAKER_06'},
{'start': 17.092, 'end': 18.323, 'speaker': 'SPEAKER_06'},
{'start': 25.968, 'end': 26.66, 'speaker': 'SPEAKER_01'}]
```

Here, the start and end markers denote the audio-timestamps corresponding to the beginning and end of the diarized segment, and the speaker denotes the speakerID corresponding to that segment. Note that there can diarization segments with multiple overlapping timestamps, if the original audio has overlapping conversation.

Language Filtering. As a next step, we identify the primary language of the audio using Whisper (Radford et al., 2023) and filter out all non-english audios.

Transcription Generation. Next, we aim to provide paired text annotations for all of the raw audios in our corpus. For this, we first used the Whisper model (Radford et al., 2023) to transcribe the raw audio from each of the diarized output chunks. However, we noticed that the Whisper model transcriptions can tend to be quite noisy and contain some hallucinations. To ensure cleaner transcriptions, we use a post-processing transcription ensembling approach called ROVER (Fiscus, 1997) used in prior works performing transcription cleaning (Jalalvand et al., 2015). We first obtain additional speech transcriptions from an internal SIRI transcription model and Nvidia-Parakeet-TDT-CTC. We then apply the ROVER post-processing method using the three candidate transcriptions from Whisper, SIRI and Parakeet. We use the ensembled transcription as our text annotations for subsequent steps. We provide some examples of the individual model-based transcriptions and the final ROVER-ensembled transcriptions below:

1026 Whisper: "And I don't think it was a compliment. Yeah." 1027 SIRI: "And I don't think it as a compliment." 1028 Parakeet: "And I don't think it's compliment yeah." 1029 ROVER-ensembled: "And I don't think it was a compliment. Yeah." 1030 1031 Whisper: "Yeah, I was just never sure if it meant like someone who was left be-1032 hind by fashion like..." 1033 1034 like" 1035

1036

1037

1039 1040

1041

1043

1044

1045 1046

1047

1048 1049

1050

1051

1052

1053

1054

1055

1056

1057 1058

1074 1075

1078 1079 SIRI: "Yeah, I was just never sure if it meant like someone who was left behind by fashion

Parakeet: "Yeah, I was just never sure if it meant like someone who was left behind by fashion like"

ROVER-ensembled: "Yeah, I was just never sure if it meant like someone who was left behind by fashion like"

Transcription Filtering. Despite the ROVER post-processing, we still find that a lot of annotations are low-quality including empty transcription texts and containing several repetitions. We filter out samples with such faulty transcriptions. For detecting repetition, we use a heuristic n-gram based approach. We first tokenize each transcription using a pretrained SentencePiece (Kudo & Richardson, 2018) tokenizer. We then search for unique 15-gram spans in the tokenized text. If we find that a 15-gram span occurs more than 5 times in the entire sequence, we discard that sample.

Interleaved Chunking. The last step in our pipeline is the interleaved chunking stage, which constructs the final audio-text chunks used for interleaved training. As described in the main text, we study two chunking strategies:

- 1. Coarse chunking. Here, we aim to have relatively long audio-text chunks. To do this, we continually merge consecutive audio segments based on the diarization outputs while they have the same speakerID. While merging the segments, we concatenate the corresponding text transcriptions of each audio segment, separated by a white-space, to yield the merged text transcription for the merged audio.
- 2. Fine chunking. Since the original diarized output segments already yield relatively short chunks, we do not do any post-processing on the output segments and directly use them as our audio-text chunks for interleaved training.

For both chunking strategies, we additionally filter out any audio-text chunks where the audio chunk is smaller than 0.2 seconds.

B DETAILS OF SYNTHETIC DATASETS

B.1 Knowledge-rich domains used for synthetic datasets

In this section, we provide a list of knowledge-rich domains we use for domain-filtering as the first step in our pipeline for constructing synthetic datasets:

- 1. https://www.numerade.com/home/
- 2. https://www.brainscape.com
- 3. https://brainly.com

- 4. https://www.chegg.com/
- 5. https://www.proprofs.com
- 6. https://www.schoolsolver.com
- 7. https://www.studypool.com
- 8. https://www.symbolab.com
- 9. https://www.justia.com
- 10. https://www.askalawyeroncall.com
- 11. https://freelawchat.com
- 12. https://www.healthtap.com
- 13. https://www.24houranswers.com
- 14. https://web2.0calc.com
- 15. https://myhomeworkapp.com
- 16. https://www.justanswer.com
- 17. https://quizlet.com/

B.2 PROMPT

Extraction prompt for *Krist.* To extract and lightly rewrite the text content from the HTML using gpt-4o-mini, we use the following prompt:

Extract the useful (non-boilerplate) text from the following HTML content into well-formatted plaintext, please. There is no need to retain hyperlinks out of the page, they can be dropped. Output the content in mark up tags as show below.

```
'''plaintext
{
  <well formatted plain text here>
}
{html_content}
```

Question validation prompt for *Quest.* To validate and filter out questions that are incorrectly formatted / extracted from the HTML, we use the following prompt to gpt-40:

 Here is a problem that you do not need to solve: {question}

Your task: Don't try to solve the problem, instead, do a brief free-form analysis, then output results for the following fields:

complete: Values choose between (you can't use any other values)

True - The problem is complete: it asks a clear and understandable question, and does not depend on any missing or unseen visual elements such as figures, graphs, tables, or images. False - The problem is incomplete: it is ambiguous, unanswerable, or relies on external content (e.g., a graph or diagram) that is not provided.

is_question: Values choose between (you can't use any other values)

True - The problem is asking a specific question (e.g., it requests the value of an expression, a numerical answer, or a specific outcome.

False - The problem is not a question (e.g., it is a statement, conversation, or unrelated content).

Question answering prompt for *Quest.* Finally, we prompt gpt-40 to answer with a chain-of-thought to each verified question using:

```
Please answer the following question. Let's think step by step. {question}
```

B.3 Examples of invalid questions in Quest

Previously, we presented the prompt used for validating and filtering out incomplete or incorrectly extracted questions. Since we score for *question completeness* and *question validity*, our filtering mechanism only keeps questions that are marked as complete *and* valid. Here, we show examples of questions that were marked as invalid i.e. marked as incomplete, invalid, or both.

```
1165
          Question
1166
          Example of mechanical?
1167
          complete: False
1168
          is_question: False
1169
1170
          Question
1171
          How does this picture show social impacts of imperialism? helppp me
1172
          complete: False
1173
          is_question: True
1174
1175
          Question
1176
          Minimum duration for diagnosis for: Selective Mutism
1177
          complete: True
          is_question: False
1178
1179
          Question
1180
          Audience analysis examples
1181
          complete: False
1182
          is_question: False
1183
1184
```

C TRAINING DATA STATISTICS

 Text-only dataset. For our text-only continued pretraining dataset, we use the dataset used in the continual pretraining experiments of Li et al. (2025a), which roughly comprises of 2.2T tokens.

Speech-text datasets. Here, we provide the exact details of all our speech-text training data sources. Note that since our tokenizer processes audio at 12.5Hz, our token yield per second is 12.5 speech tokens. Hence, an hour of audio (3600s) corresponds to 45k speech tokens. In Tab. 7, for each dataset, we report the number of raw hours of speech content along with the total number of speech tokens. As evident, web-crawl data contains the most number of unique tokens followed by Krist and Quest.

Table 7: Training Data Statistics.

Training dataset	# Hours	# Speech tokens
Web-crawl	8.03M	361.3B
Krist	4.72M	212.4B
Quest	0.86M	38 B

C.1 DETAILS OF DATA MIXTURES FOR SYNTHETIC DATA EXPERIMENTS

Here, we break down the exact token counts used for each data mixture in the experiments in Tab. 2. Remember that we train for a total of 200k steps with a batch-size of 512 and sequence-length of 16, 384 yielding 1.67T multimodal tokens for the full training run. For each experiment, we use 60% text-only and 40% speech-text mixing ratio. Hence, the text-only ratio corresponds to \sim 1T tokens. The speech-text ratio corresponds to the remaining \sim 670B tokens. Now, in Tab. 8, we report for each data source (text-only, web-crawl, Krist and Quest), the exact mixing proportion in the training mixture (%mix), total number of tokens in the training mixture (#toks) and the number of repeats (epochs) of the original data source (#repeats) used across all our experiments in Tab. 2. As is evident from the table, due to the heterogenity of data sources and their corresponding token-sizes, it is quite complex to determine an optimal mixing proportion.

Table 8: Data mixture statistics for experiments in Tab. 2.

Training dataset	Text-only dataset		Web-crawl		Krist			Quest				
Trumming durantee	%mix	#toks	#repeats	%mix	#toks	#repeats	%mix	#toks	#repeats	%mix	#toks	#repeats
Web-crawl 100%	0.60	1T	0.45	0.40	670B	1.85	0.00	0.00	0.00	0.00	0.00	0.00
Web-crawl 53% + Krist 47%	0.60	1T	0.45	0.21	355B	0.98	0.19	315B	1.48	0.00	0.00	0.00
Web-crawl 66% + Quest 34%	0.60	1T	0.45	0.26	442B	1.22	0.00	0.00	0.00	0.14	228	6.00
Web-crawl 59% + Quest 6% + Krist 35%	0.60	1T	0.45	0.24	395B	1.09	0.14	232B	1.10	0.02	43B	1.13
Web-crawl 40% + Quest 27% + Krist 33%	0.60	1T	0.45	0.16	267B	0.74	0.13	221B	1.04	0.11	182	4.79

D TRAINING DETAILS

All our models are 3.8B parameter transformer-based speech-language models. We use a global-batch-size of 512 for all our experiments. Our models use a packed-sequence-length of 16,384 tokens. We train for 200k steps in total, yielding a total of 1.67T multimodal tokens for our training runs. Using the standard 6ND rule (Kaplan et al., 2020), this equates to about 3.81×10^{22} FLOPs (note that this estimate is a rough lower bound since we do not count the FLOPs associated with the speech tokenizer in this estimate). We only tune the language model weights while keep the speech tokenizer frozen. We use a cosine-decay learning rate schedule with 1000 steps of linear-warmup. We use the AdamW (Loshchilov & Hutter, 2017) optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.95$, a peak learning rate of 3e-4, weight decay of 1e-5 and clip gradients to a max norm of 1.0. We use the axlearn (Lee et al., 2025) codebase for all our experiments using jax (Bradbury et al., 2021) and pygrain (Ritter et al., 2023) for dataloading. One training run takes approximately 7 days on 512 TPU-v6e chips.

E EXTENDED RELATED WORK

 In the main paper, we briefly described some related work in speech-language pretraining. Further, we focused on situating our work in the SpeechLM literature and emphasized the lack of data-centric research in speech-language pretraining. Here, we provide a deeper dive into SpeechLMs and reference some related data-centric work that does exist in the speech-language domain.

Speech Language Models. There has been a recent push for training end-to-end SpeechLMs (Arora et al., 2025). Early efforts like Whisper (Radford et al., 2023), SALMONN (Tang et al., 2023), and LTU-AS (Gong et al., 2023) employed multi-task pretraining to enable tasks like automatic speech recognition, emotion classification etc. Scaling these principles by increasing model-size and training compute (Chu et al., 2023; 2024; Liu et al., 2025; Geng et al., 2025; Kong et al., 2024; Ghosh et al., 2025; Goel et al., 2025) has yielded continued gains. Further works considered pretraining models with speech understanding and generation capabilities (Lakhotia et al., 2021; Algayres et al., 2023; Hassid et al., 2023; Nguyen et al., 2025b; Nachmani et al., 2023; Rubenstein et al., 2023; Zhang et al., 2023; Défossez et al., 2024). More recently, models like Kimi-Audio (Ding et al., 2025), Step-Audio-2 (Wu et al., 2025), Baichuan-Audio (Li et al., 2025b), GLM-4-Voice (Zeng et al., 2024a), and MiMo-Audio (Xiaomi, 2025) have emerged as strong foundation models that seamlessly perform several tasks, including spoken-question answering. While demonstrating impressive performance, details behind their data curation strategies are scant. Through our controlled experiments, we aim to fill this gap by shedding light on how to effectively construct speech-text pretraining datasets.

Data Curation for Speech-Language Models. Whisper (Radford et al., 2023) was one of the first works to effectively leverage web-scale data for training a multi-task speech-text model, using a dataset of 680k hours. Attempting to openly reproduce the original Whisper dataset, (Ngo et al., 2025) introduced OLMOASR-POOL, a dataset of 3M hours of audio and 17M transcripts. They conducted heuristic-based filtering on their data pool, showcasing benefits on ASR tasks. Tian et al. (2024) and Peng et al. (2025) similarly conducted comprehensive studies to understand the effects of data heterogenity, ASR error rate based filtering and LLM-based transcription rephrasing, while training Whisper-style models. However, these efforts were limited to training models that were primarily capable of performing ASR tasks. The data curation literature in the end-to-end SpeechLM literature is much more sparse. Kimi-Audio (Ding et al., 2025) describes their speech-text dataset construction pipeline, beginning from 13M audio hours and processing them into speech-text interleaved training data. However, why certain design decisions were taken remain unanswered. Contrarily, Zeng et al. (2024b) constructed synthetic interleaved data sourced from high-quality text pretraining data, but yet again omit clear details on key design choices. MiMo-Audio (Xiaomi, 2025) scaled up their training dataset size by an order of magnitude to an unprecedented 100M hours of audio data. While they showcased the benefits of dataset quantity using few-shot experiments, they did not conduct any explicit controlled experiments to justify the filtering and curation decisions they made. In our work, we aim to fill this gap on the data-centric side of SpeechLMs, by describing and understanding data curation pipelines for speech-text interleaved pretraining through three key questions around interleaved data chunking, synthetic dataset construction and chunk sampling schemes during interleaved training.

F DETAILS AND EXAMPLES OF SQA EVALUATION DATASETS

We aim to evaluate the *speech-to-text transfer* capability of SpeechLMs, where the model is asked a question in speech and tasked with responding in text ($S \rightarrow T$). In the literature, there is a lack of standardized evaluations for this task of Spoken-Question-Answering (SQA). While efforts like Spectron-LM (Nachmani et al., 2023) and Voxtral (Liu et al., 2025) have open-sourced some evaluation sets, they use different text-to-speech engines and generation parameters for synthesizing the spoken questions, rendering comparisons across different models unfair. Moreover, these datasets only consist of a question and answer, requiring models to generate free-form text outputs. However, prior works in LM evaluation standardization (Gu et al., 2024; Allal et al., 2025; Li et al., 2024; Brown et al., 2020) recommend using a cloze-form of MCQ evaluation for evaluating base-models with question-conditioned completion log-probabilities rather than decoding free-form text outputs. The log-probability method removes evaluation confounds such as decoding temperature, sampling method and other decoding parameters, which are known to induce large variance (Hochlehnert et al., 2025). Therefore, we construct a standardized SQA evaluation suite of three datasets—Spoken-LLaMA-Questions, Spoken-Web-Questions and Spoken-TriviaQA. We source the raw audio questions from OpenAudioBench (Li et al., 2025b). We then prompt qpt-40-mini with the original text question and answer of each sample to provide a set of three distractor choices (the prompts for generating choices are in Appx. G). Hence, our final evaluation datasets consist of a spoken-question and 4 choices, with one correct answer (chance-level is 25%). In Tab. 9, we provide details about the number of test samples, the TTS engine used for synthesizing the speech questions, and the links to the original audio source files.

Table 9: Details of SQA evaluation datasets.

Evaluation Dataset	Num. samples	Chance%	TTS Engine	Audio Source
Spoken-LLaMA-Questions	300	25%	Google Cloud TTS	Link
Spoken-TriviaQA	1000	25%	Baichuan-Audio TT	S Link
Spoken-Web-Questions	1000	25%	Baichuan-Audio TT	S Link

Below, we also provide a few examples from each evaluation dataset, with the question (in text), choices, and the ground-truth answer.

Spoken-LLaMA-Questions

Question: What is the capital of France?
Choices: Paris, London, Berlin, Madrid
Ground-Truth: Paris

Question: Which river is the longest in South America?
Choices: Nile, Amazon, Paraná, Orinoco
Ground-Truth: Amazon

Spoken-TriviaQA

Question: Who was Jackie Kennedy's second husband?
Choices: John F. Kennedy, Robert F. Kennedy, Frank Sinatra, Aristotle Onassis
Ground-Truth: Aristotle Onassis

Question: What is the oldest vegetable known to man?
Choices: Carrot, Potato, Pea, Onion
Ground-Truth: Pea

Spoken-Web-Questions

Question: What language do most Italians speak?
Choices: Italian, French, Spanish, German
Ground-Truth: Italian

Question: Who did Shaq first play for?
Choices: Los Angeles Lakers, Miami Heat, Boston Celtics, Orlando Magic
Ground-Truth: Orlando Magic

Evaluation details. We use log-likelihood based scoring for our evaluation protocol following standard language modeling works (Brown et al., 2020; Allal et al., 2025; Gu et al., 2024).

For each test sample and each answer-choice (out of 4 total choices), we use the following cloze-form to prompt the model:

Question:\n<question-in-audio>\nAnswer:<answer-choice>

Then, we compute the completion log-probability for each of the 4 answer choices. We normalize the completion log-probability by answer length to prevent biasing against long answer choices. A question is marked correct if the model assigns highest normalized log-probability to the ground-truth answer. We use standard accuracy metric (random chance level is 25%) for reporting results.

PROMPTS FOR GENERATING DISTRACTOR CHOICES FOR EVALUATION SETS We use the following prompt for generating the distractor options for Spoken-LLaMA-Questions and Spoken-TriviaQA. SYSTEM PROMPT You are a helpful assistant. INPUT PROMPT I will give you a simple question and answer pair. This pair comes from an evaluation dataset. I am trying to convert it into an MCQ format dataset. You have to give three more plausible distractor options that I can use along with the correct option to create the MCQ test set. Give the three distractor options one after the other, comma-separated, all in one line. Here are a few examples: Input: Question: What colour is the sky? Answer: blue Output: green,red,yellow Input: Question: What season comes after spring? Answer: summer Output: winter, monsoon, autumn I will give you the question and the answer now. Remember, please give the three options in one line, comma-separated. Question: <question> Answer: <answer>

For *Spoken-Web-Questions*, as there can be multiple correct answers for a question, we pick the first reference answer as ground-truth and use the following prompt for generating distractor options.

1515 SYSTEM PROMPT

You are a helpful assistant.

INPUT PROMPT

I will give you a simple question and answer pair. This pair comes from an evaluation dataset. Note that the answer might be one of out many possible correct answers. I am trying to convert it into an MCQ format dataset. You have to give three more plausible distractor options that I can use along with the correct option to create the MCQ test set. Since the provided answer might be one of many possible correct answers, ensure that the distractor options you provide are definitely incorrect for the given question. For example, if the question is "What is a leap year?" and the answer I provide is 2004, do not give distractor options like 2000 or 2012. Give the three distractor options one after the other, comma-separated, all in one line.

Here are a few examples:

Input:

Question: What colour is the sky?

Answer: blue

Output:

green,red,yellow

Input:

Question: What season comes after spring?

Answer: summer

Output:

winter, monsoon, autumn

I will give you the question and the answer now. Remember, please give the three options in one line, comma-separated.

Question: <question>
Answer: <answer>

H DIVERGENCE ANALYSIS BETWEEN MODALITY DISTRIBUTIONS

In this section, we describe in detail the exact setup used for our analysis in Sec. 4.1.

We start with a spoken question-answering test set. Each test sample consists of (q_a, q_t, gt) triplets, where q_a denotes the spoken question in audio modality, q_t denotes the question in text modality, and gt denotes the ground-truth answer in text modality.

Goal. We aim to measure the divergence between the token-wise teacher-forced (Williams & Zipser, 1989) conditional probability distributions of the audio and text modality. That is, we compare the next-token distributions under audio vs. text question conditioning, evaluated along the same ground-truth (GT) answer path (the answer is always in text modality).

Notation. For each test sample s, let $\{t_1, t_2 \cdots t_m\}$ and $\{a_1, a_2 \cdots a_n\}$ represent the question tokens in text and audio modality respectively. That is, the tokenized representation of q_t is $\{t_1, t_2 \cdots t_m\}$ and the tokenized representation of q_a is $\{a_1, a_2 \cdots a_n\}$. For brevity, let us denote these tokenized representations as $t_{1:m}$ and $a_{1:n}$. Note that since the length of the question tokens in text and audio modalities might differ, it is possible that $n \neq m$. Let $\{g_1, g_2 \cdots g_o\}$ represent the ground-truth answer tokens in text modality i.e. the tokenized representation of gt is $\{g_1, g_2 \cdots g_o\}$. Again, for brevity, we denote this as $g_{1:o}$. Let V be the vocabulary of the SpeechLM.

For a given test sample s, for each answer token $i \in \{1, 2 \cdots o\}$, we define the teacher-forced next-token distributions as:

$$P_{\text{aud},i}^{(s)}(v) = \Pr_{\theta} (X = v \mid a_{1:n}, g_{1:i-1}), \quad v \in V,$$
(1)

$$P_{\text{text},i}^{(s)}(v) = \Pr_{\theta} (X = v \mid t_{1:m}, g_{1:i-1}), \quad v \in V.$$
 (2)

where $\Pr_{\theta}(X=v|Y)$ represents the conditional probability distribution for all values $v \in V$, conditioned on the previous context Y.

Per–token divergences. We now compute (1) forward KL, (2) reverse KL, and (3) Jensen–Shannon (JS) divergence at each step i, between the two next-token distributions:

$$D_{\text{KL}\to}^{(s)}(i) = \sum_{v \in V} P_{\text{aud},i}^{(s)}(v) \log \frac{P_{\text{aud},i}^{(s)}(v)}{P_{\text{text},i}^{(s)}(v)},\tag{3}$$

$$D_{\text{KL}\leftarrow}^{(s)}(i) = \sum_{v \in V} P_{\text{text},i}^{(s)}(v) \log \frac{P_{\text{text},i}^{(s)}(v)}{P_{\text{and }i}^{(s)}(v)}, \tag{4}$$

$$D_{\rm JS}^{(s)}(i) = \frac{1}{2} D_{\rm KL} \left(P_{{\rm aud},i}^{(s)} \parallel M_i^{(s)} \right) + \frac{1}{2} D_{\rm KL} \left(P_{{\rm text},i}^{(s)} \parallel M_i^{(s)} \right), M_i^{(s)} = \frac{1}{2} \left(P_{{\rm aud},i}^{(s)} + P_{{\rm text},i}^{(s)} \right). \tag{5}$$

Answer–span aggregation (per example). To get a mean divergence value per sample, we average the per-token divergences over the answer length o (masking any padded positions in practice):

$$\overline{D}_{\text{KL}\to}^{(s)} = \frac{1}{o} \sum_{i=1}^{o} D_{\text{KL}\to}^{(s)}(i), \qquad \overline{D}_{\text{KL}\leftarrow}^{(s)} = \frac{1}{o} \sum_{i=1}^{o} D_{\text{KL}\leftarrow}^{(s)}(i), \qquad \overline{D}_{\text{JS}}^{(s)} = \frac{1}{o} \sum_{i=1}^{o} D_{\text{JS}}^{(s)}(i). \tag{6}$$

The distribution of these per-sample mean divergences is what we plot in Fig. 5 and Appx. H.1.

Dataset-level metrics. Over each test set S we also report the dataset means across metrics in Tab. 10:

$$\mathcal{D}_{\mathrm{KL}\to} = \frac{1}{|S|} \sum_{s \in S} \overline{D}_{\mathrm{KL}\to}^{(s)}, \qquad \mathcal{D}_{\mathrm{KL}\leftarrow} = \frac{1}{|S|} \sum_{s \in S} \overline{D}_{\mathrm{KL}\leftarrow}^{(s)}, \qquad \mathcal{D}_{\mathrm{JS}} = \frac{1}{|S|} \sum_{s \in S} \overline{D}_{\mathrm{JS}}^{(s)}. \tag{7}$$

H.1 More results across different metrics and test sets

In the main paper Sec. 4.1, we showcased the divergence plots between the conditional next-token distributions, on the Spoken-LLaMA-Questions test with the reverse KL-divergence metric only. Here, we showcase the divergence distributions across all three of our test sets—Spoken-LLaMA-Questions, Spoken-Web-Questions and Spoken-TriviaQA—across three divergence metrics—Forward KL Divergence, Reverse KL Divergence and Jensen Shannon Divergence. The plots for Spoken-LLaMA-Questions are in Fig. 9, for Spoken-Web-Questions are in Fig. 10, and for Spoken-TriviaQA are in Fig. 11. Furthermore, in Tab. 10, we report the mean values of the divergence distributions obtained. Across all plots and the table, we observe that our data interventions consistently close the distribution mismatch between the conditional probability distributions of audio and text modalities. This suggests that our data intervention implicitly induce a self-distillation behaviour (Zhang et al., 2021; Mobahi et al., 2020; Zhang et al., 2019) in our trained SpeechLMs. Such an implicit "distillation through data" property has also been observed in prior works in the multimodal and language domains (Udandarao et al., 2025; Rawat et al., 2024; Wang et al., 2024; Sachdeva & McAuley, 2023; Wang et al., 2018). Further, Wang et al. (2025a) showed that explicitly applying a cross-modal distillation objective further helps to reduce the modality distribution gap, and our results further implicitly confirm this.

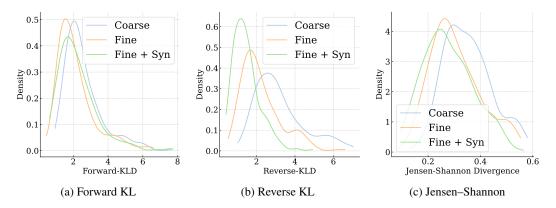


Figure 9: Conditional-distribution divergences on Spoken-LLaMA-Questions.

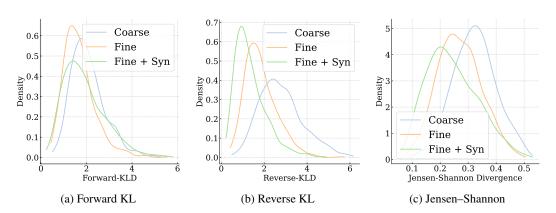


Figure 10: Conditional-distribution divergences on Spoken-Web-Questions.

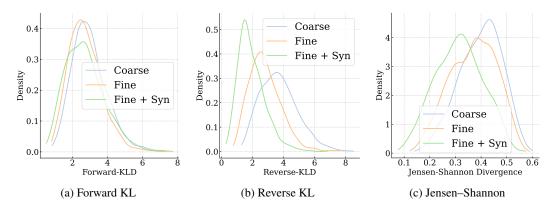


Figure 11: Conditional-distribution divergences on Spoken-TriviaQA.

Table 10: **Dataset-level means of all divergence metrics b/w conditional next-token distributions.** We report the means of all three divergence distributions (as computed in eq. (7)). *FKL* represents forward KL-divergence, *RKL* is reverse KL-divergence and *JSD* is Jensen-Shannon divergence.

Method	Spoken-Web-Questions			Spoken-TriviaQA			Spoken-LLaMA-Questions			
	\overline{FKL}	RKL	JSD	\overline{FKL}	RKL	JSD	\overline{FKL}	RKL	JSD	
Coarse	2.07	2.78	0.32	2.97	3.70	0.40	2.57	3.20	0.35	
Fine	1.68	1.84	0.27	2.72	2.80	0.36	2.15	2.21	0.30	
Fine + Syn	1.90	1.35	0.24	2.71	1.94	0.31	2.23	1.47	0.27	

I TOPIC DOMAIN ANALYSIS

I.1 DETAILS ABOUT TOPIC DOMAIN CLASSIFIER

For conducting the topic domain analysis in Fig. 6, we used the topic domain classifier that was released by (Wettig et al., 2025). The classifier is a gte-base-en-v1.5 that was fine-tuned on web-texts annotated by LLaMA models. We used the No-URL version of the classifier that takes only the raw text as input and classifies it into one of 24 output classes. For getting the topic distribution of each of our datasets, we randomly sample 5000 examples, concatenate all the text chunks from each example (for web-crawled data, these are the annotated transcriptions while for synthetic data, these are the source text data samples), and use that as input to the topic classifier.

I.2 TOPIC DISTRIBUTION FOR SPOKEN-WEB-QUESTIONS

In Fig. 12, we showcase the topic distribution of Spoken-Web-Questions. Similar to the takeaways in Fig. 6, we find that some of the topics that Spoken-Web-Questions contains are severely under-represented in the web-crawled dataset while being represented adequately in the synthetic datasets. This further corroborates our findings that synthetic datasets help close the distribution mismatch between the web-crawled dataset and the evaluation datasets. Our findings regarding the under-representation of concepts in web-crawled datasets have also been echoed in the language and vision domains (Wiedemer et al., 2025; Parashar et al., 2024; Elazar et al., 2023; Kandpal et al., 2023; Udandarao et al., 2024; Zhao et al., 2024; Samuel et al., 2024; Dodge et al., 2021).

I.3 A MORE FINE-GRAINED TOPIC DISTRIBUTION ANALYSIS

For all the topic domain analyses we have conducted previously, we used a coarse-level topic classifier that could categorize between 24 different topics. Here, we use a more fine-grained topic classifier that can produce a finer-grained categorization into 67 different topics. We use the finefineweb-domain-fasttext-classifier, which is a bi-gram fasttext model that was used for curating the FineFineWeb dataset (Zhang et al., 2024a). We use the same procedure as

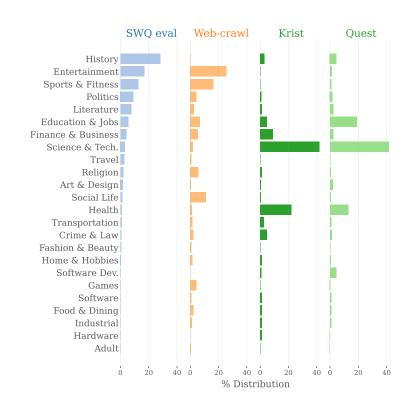


Figure 12: Topic domain distribution for Spoken-Web-Questions eval and training datasets.

before for annotating our evaluation and training datasets. We plot the fine-grained topic distributions for *Spoken-LLaMA-Questions* in Fig. 13, *Spoken-TriviaQA* in Fig. 14 and *Spoken-Web-Questions* in Fig. 15, along with all training datasets. Across all the plots, our findings from Figs. 6 and 12 hold—our synthetic datasets increase the diversity and topic coverage of our training data distribution, thereby more closely matching the distribution of concepts encompassed in the evaluation datasets. This helps improve model generalization, yielding better downstream performance.

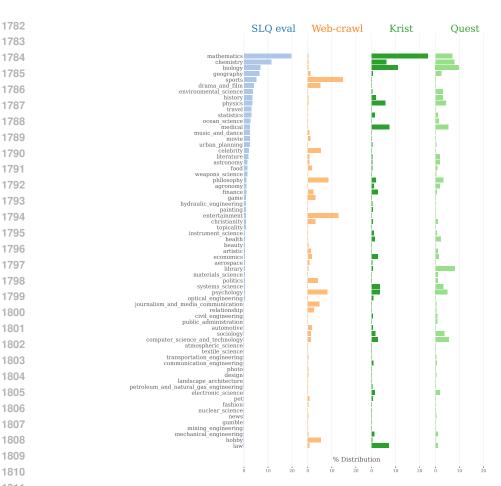


Figure 13: Fine-grained topic domain distribution for Spoken-LLaMA-Questions eval and training datasets.

DETAILS ABOUT CONTAMINATION ANALYSIS

J.1 **EXAMPLES OF CONTAMINATED MATCHES**

In this section, we show some examples of the matches we get from our contamination identification procedure. For each match, we show the training dataset, the training sample, the contaminated test sample, the test dataset it belongs to, and the contaminated n-gram span.

```
Train dataset: Quest
Train sample:
```

What is the definition of vitreous? The word derives from Latin viteus, "of glass," and is used to describe either a glass-like quality or the glass-like substance filling the eye. Vitreous (adjective): 1. Having the appearance or properties of glass; glassy, transparent, brittle. 2. In anatomy, relating to the vitreous humor or vitreous body—the clear, gelatinous substance filling the space between the lens and the retina of the eye.

```
Test dataset: Spoken-TriviaQA
```

Test sample:

1807

1811

1812

1813 1814 1815

1816 1817

1818 1819

1820

1821

1822 1823

1824

1825

1826

1827

1828

1829

1830

1831

1832

1833

1834

1835

What is the thick watery substance filling the space between the lens and the retina of the eye?

Contaminated span:

substance filling the space between the lens and the retina of the eye



Figure 14: Fine-grained topic domain distribution for *Spoken-TriviaQA* eval and training datasets.

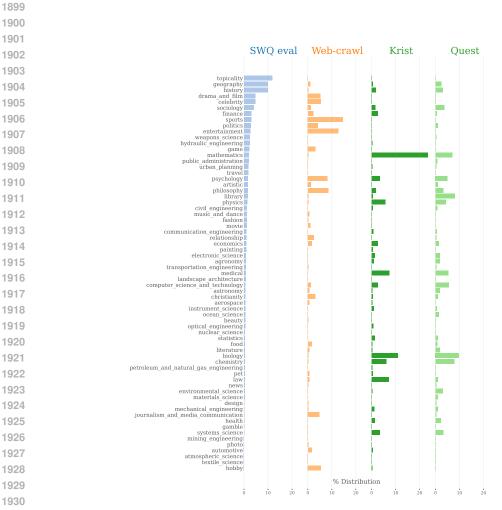


Figure 15: Fine-grained topic domain distribution for *Spoken-Web-Questions* eval and training datasets.

1944 1945 Train dataset: Quest

Train sample:

1946

1947

1948

1949

1950

1953

1954

1955

1956

1961

1962

1963

1964

1965

1966

1968

1969

1970

1971 1972

1974

1975

1976

1978

1981

1982

1984

1985 1986 1987

1988 1989

1992

1996

When did Arthur & Catherine marry? Prince Arthur, the eldest son of Henry VII, married Catherine of Aragon — daughter of Ferdinand II of Aragon and Isabella I of Castile. Their wedding took place at St. Paul's Cathedral in London on 14 November 1501. Arthur and Catherine were married on 14 November 1501.

Test dataset: Spoken-TriviaQA

Test sample:

What was founded by Ferdinand II of Aragon and Isabella I of Castile to keep Catholic orthodoxy as the major religion of their kingdoms?

Contaminated span:

ferdinand ii of aragon and isabella i of castile

1957 1958 Train dataset: Krist

Train sample:

What conclusions can be drawn about the USA's actions in the 1920s and 1930s? One conclusion to this statement, which seems to be addressing the approach to foreign policy during the period, might be "...reflected a strong, if uneven, commitment to isolationism." On the one hand, the United States was fairly steadfast in its unwillingness to get directly involved in the affairs of the world, particularly Europe. Except for a few non-binding pacts and negotiations over the repayment of reparations and war debts, the United States remained generally aloof from European affairs during the 1920s.

Test dataset: Spoken-TriviaQA

1967 Test sample:

What was the name of the democratic government of Germany in the 1920s and early 1930s, destroyed by Adolf Hitler?

Contaminated span: in the 1920s and early 1930s,

Train dataset: Krist Train sample:

In 1912, Lenin, then in exile in Switzerland, appointed Joseph Stalin to serve on the first Central Committee of the Bolshevik Party. Three years later, in November 1917, the Bolsheviks seized power in Russia. The Soviet Union was founded in 1922, with Lenin as its first leader. During these years, Stalin had continued to move up the party ladder, and in 1922 he became secretary general of the Central Committee of the Communist Party, a role that enabled him to appoint his allies to government jobs and grow a base of political support.

Test dataset: Spoken-Web-Questions

Test sample:

what led to stalin rise in power?

Contaminated span:

to serve on the first central committee of the bolshevik party

J.2 Proportion of Contamination in eval datasets

Table 11: **Proportion of contamination.** For each evaluation dataset, we report the proportion of test samples detected as contaminated. We also report the absolute number of matches in brackets.

Evaluation dataset	% Contamination [# samples]					
2,44444044	Krist	Quest	All			
Spoken-Web-Questions	0.4% [4]	0.1% [1]	0.4% [4]			
Spoken-TriviaQA	2.2% [22]	0.8% [8]	2.5% [25]			
Spoken-LLaMA-Questions	6.7% [20]	0.2% [5]	7.7% [23]			

J.3 EXPANDED DESCRIPTION OF SIGNIFICANCE TESTING SETUP AND RESULTS

Null hypothesis. We start from the full test set (containing contaminated samples). In our significance test, we test *whether removing contaminated test items reduces accuracy beyond what would be expected under random removal of an equal number of items.* Formally, for accuracy A, the null is:

$$H_0: A_{\text{clean}} \sim \text{distribution of } A_{\text{rand}},$$

i.e., the clean accuracy is not lower than the random-removal distribution. Because the contamination claim is directional (contamination would inflate accuracy), we use a *one-sided* test.

Test procedure. For each training mix and dataset from Sec. 3.4, we compute: (i) *Full* accuracy on the full test set; (ii) *Clean* accuracy after removing all known contaminated items; (iii) a *random-removal baseline* by drawing 100 random subsets (without replacement) of the same size as the contaminated set, recomputing accuracy on the remaining items each time. Accuracies for (ii) and (iii) are computed over the reduced denominators (remaining items). From the bootstrap distribution we report the mean and 95% percentile CI and compute the empirical one-sided *p*-value as:

$$p = \Pr(A_{\text{rand}} \leq A_{\text{clean}}),$$

This p-value is appropriate for the hypothesis that contamination inflates accuracy (so clean should be lower if inflation is present). With 100 replicates, the p-value granularity is 0.01. Hence, we report p<0.01 when no replicate from the bootstrap distribution is as low as the clean accuracy.

Results and interpretation. Tables 12–14 summarize results for Spoken-TriviaQA, Spoken-LLaMA-Questions, and Spoken-Web-Questions. We highlight the difference $\Delta = \text{Clean}-\text{RandMean}$ and give the decision at a significance level $\alpha = 0.01$.

Table 12: One-sided contamination test on STQ (N=1000).

Data mix	Full (%)	Clean (%)	Random mean (95% CI) (%)	Δ (pp)	One-sided p	Decision
Web-crawl 53% + Krist 47%	29.20	29.03	29.22 [28.77, 29.64]	-0.19	0.32	Fail to reject H_0
Web-crawl 66% + Quest 34%	34.70	34.56	34.73 [34.26, 35.28]	-0.17	0.38	Fail to reject H_0
Web-crawl 59% + Quest 6% + Krist 35%	30.80	30.46	30.81 [30.36, 31.18]	-0.35	0.09	Fail to reject H_0
Web-crawl 40% + Quest 27% + Krist 33%	31.70	31.59	31.70 [31.28, 32.10]	-0.11	0.41	Fail to reject H_0

Table 13: One-sided contamination test on SLQ (N=300).

Training mix	I	Full (%)	1	Clean (%)	ī	Random mean (95% CI) (%)	ī	Δ (pp)	One-sided p	1	Decision
Web-crawl 53% + Krist 47%	†	52.00	Ť	50.54	i	52.16 [50.54, 53.62]	İ	-1.62	0.10	Ť	Fail to reject H_0
Web-crawl 66% + Quest 34%		66.33		66.79		66.34 [64.62, 68.06]		+0.45	0.82		Fail to reject H_0
Web-crawl 59% + Quest 6% + Krist 35%		50.33	ı	48.01	İ	50.33 [48.91, 51.62]	İ	-2.32	< 0.01	ı	Reject H_0
Web-crawl 40% + Quest 27% + Krist 33%		49.33	İ	47.29	İ	49.43 [47.65, 50.90]	İ	-2.14	0.02	ı	Fail to reject H_0

Table 14: One-sided contamination test on SWQ (N=1000).

Training mix	F	Full (%)	C	Clean (%)	Random mean (95% CI) (%)	Δ (pp)	One-sided p	Decision
Web-crawl 53% + Krist 47%		43.40	1	43.27	43.39 [43.22, 43.57]	-0.12	2	0.23	Fail to reject H_0
Web-crawl 66% + Quest 34%		42.70	ı	42.57	42.69 [42.47, 42.87]	-0.12	2	0.20	Fail to reject H_0
Web-crawl 59% + Quest 6% + Krist 35%		43.80		43.67	43.80 [43.62, 43.98]	-0.13	3	0.19	Fail to reject H_0
Web-crawl 40% + Quest 27% + Krist 33%		43.30		43.17	43.29 [43.07, 43.47]	-0.12	2	0.23	Fail to reject H_0

Takeaways. Across STQ and SWQ, clean accuracies consistently fall within the random-removal confidence intervals. Therefore, we find no significant contamination-driven inflation. For SLQ, the Web-crawl 59% + Quest 6% + Krist 35% mix shows a drop in clean accuracy relative to the random baseline that is statistically significant under our one-sided p-test (p<0.01), consistent with contamination inflating test performance. However, for the other three data mixes we again see no significant evidence of inflation, under our testing setup. Hence, overall we conclude that contamination does not have a major effect on inflating model performance.

J.4 LIMITATIONS OF OUR CONTAMINATION ANALYSIS

Post-hoc analysis. Our contamination analysis is entirely post-hoc, after training of a model is complete. In the ideal case, one would decontaminate the training sets with respect to the test sets

a-priori (Beyer et al., 2024; Zhai et al., 2022; Oquab et al., 2023; Trinh & Le, 2018; Gao et al., 2020; Mizrahi et al., 2025; Allal et al., 2025; OLMo et al., 2024). In practice, however, this is unrealistic, since this assumes prior knowledge of all possible test sets that the model may encounter in the wild. Infact, several popular language model trainers do not decontaminate their training sets precisely for this reason (Su et al., 2024; Weber et al., 2024; Maini et al., 2025; Rae et al., 2021; Penedo et al., 2023; Kandpal et al., 2025). Further, while we acknowledge that our post-hoc contamination analysis can be limiting and would benefit from a more causal treatment such as in works like (Li et al., 2024; Soldaini et al., 2024; Bordt et al., 2024; Jiang et al., 2024), we however note that the downside of such a causal analysis is the significant overhead of re-training our models. Hence, we also note that many works in the literature refrain from a fully causal treatment of contamination (Radford et al., 2019; Brown et al., 2020; Dubey et al., 2024; Achiam et al., 2023).

Language-only detection. Our contamination detection only operates on the seed text-datasets that we generate our synthetic datasets from. We have not done any contamination analysis between the spoken question audios in our test sets with the audios in our training sets (we note that prior works in speech-language processing also mainly do contamination analysis at the text-level (Ngo et al., 2025; Tseng et al., 2025)). While this is a reasonable proxy for our synthetic datasets, such a method might not transfer well for decontamination analyses of web-crawled datasets. This is because many of the speech transcriptions of the web-crawled speech might be noisy, incorrect or contain hallucinations induced by the transcription model. Hence, measuring, detecting and quantifying contamination on the audio modality is an important research problem that warrants futher attention.

Testing on non-contaminable benchmarks. While research in optimal ways to do test-set contamination in language models is still nascent, many works take the alternate approach of building benchmarks that are by construction non-contaminated (Ghosh et al., 2024; White et al., 2024; Zeng et al., 2025; Wildman et al., 2025; Jain et al., 2024; Zhang et al., 2024b; Yang et al., 2023; Srivastava et al., 2024). We note that there is a huge gap in such robust evaluations in the speech-language modeling community, and striving for better benchmarks would enable stronger significance in results, while diminishing the impacts of train-test contamination on downstream model performance.

K LIMITATIONS AND DISCUSSION

While we conducted extensive experiments to study the three data-centric questions we outlined in Fig. 1, there are still a few limitations in our work that can be improved upon:

Model sizes. All our experiments were at the 3.8B parameter scale. While our results are strong (outperforming models that are $3 \times$ the size), it would still be interesting to explore if our data-centric strategies would hold at larger model scales. While papers like DataComp-LM (Li et al., 2024) and DataComp-CLIP (Gadre et al., 2023) suggest transferability of data curation methods across model scales, recent work in language modeling has posited that there may be trade-offs when applying data curation across different model sizes (Mizrahi et al., 2025). To the best of our knowledge, no existing work showcases such trade-offs in the SpeechLM community. It would be an interesting direction to explore the interaction of data recipes and model scale.

More speech-text tasks. Since the focus of our work was mainly on improving spoken question-answering capabilities of SpeechLMs, all our experiments used the standard benchmarks that are prevalent in the literature for our task of interest (Liu et al., 2025; Xiaomi, 2025; Li et al., 2025b; Ding et al., 2025). We therefore did not explore how our models would perform on more targeted tasks like automatic speech recognition, emotion recognition or text-to-speech synthesis. One caveat preventing us from a direct comparison on such tasks is that we do not employ any task-specific training, unlike other SpeechLMs that explicitly add in a task-specific component into their training mixture (e.g., ASR-specific training datasets).

Transfer to post-training? In our work, we only explored speech-language pretraining. Hence, all our data-centric methods were only tested for the pretraining phase. Therefore, our model checkpoints are inherently base-models, which cannot be used in an assistant manner. It is an interesting open question whether our data intervention strategies would also lead to better post-training results. To the best of our knowledge, such a connection between pretraining and post-training is underexplored even in the language modeling community, with only a few works hinting at such a connection with no clear experimentation (Mizrahi et al., 2025; Maini et al., 2025).

Training from scratch. All our training runs initialize the language model backbone for our SpeechLM using a pretrained base-LM. This is the standard recipe used by almost all the existing foundation SpeechLMs (Li et al., 2025b; Défossez et al., 2024; Liu et al., 2025; Wu et al., 2025; Xiaomi, 2025; Chu et al., 2024; Ding et al., 2025; Zeng et al., 2024a). However, recent work in the vision-language literature has advocated for full native multimodal pretraining from scratch (Shukor et al., 2025), where both the language model and the modality-specific encoder/tokenizer are trained from scratch. It would be interesting to explore if our data-centric methods also enable more efficient SpeechLM pretraining from scratch in the future.

Better training recipes. In all our experiments, we freeze the speech tokenizer while only training the language model. In the SpeechLM literature, there is no strong consensus regarding freezing or unfreezing the speech tokenizer. A potential next step could be to unfreeze the tokenizer and study the transferability of our data-centric recipes. Additionally, we conduct only one continued-pretraining stage—however, recent SpeechLM works have explored more sophisticated multi-stage pipelines involving pretraining and mid-training (Wu et al., 2025; Xiaomi, 2025; Li et al., 2025b; Goel et al., 2025). It would again be interesting to test our methods in a multi-stage pipeline.