

Improving Spoken Semantic Parsing using Unpaired Text from Textual Corpora and Large Language Model Prompting

Anonymous ACL submission

Abstract

Spoken semantic parsing (SSP) involves generating machine-comprehensible parses from input speech. Training robust models for existing application domains represented in training data or extending to new domains requires corresponding triplets of speech-transcript-semantic parse data, which is expensive to obtain. In this paper, we address this challenge by examining methods that can use or generate transcript-semantic parse data (unpaired text) without corresponding speech. First, when unpaired text is drawn from existing textual corpora, we compare Joint Audio Text (JAT) and Text-to-Speech (TTS) as ways to use unpaired text to generate speech representations. Experiments on the STOP dataset show that unpaired text from existing and new domains improves performance by 2% and 30% in absolute Exact Match (EM) respectively.

Second, when unpaired text is not available from existing textual corpora, Large Language Models (LLMs) can be prompted to generate unpaired text for existing and new domains, and JAT or TTS can be used with the generated unpaired text to improve SSP. Prior work has mostly focused on using LLMs to generate synthetic data for classification tasks. In this paper, we introduce multiple prompting strategies to obtain synthetic data in existing and new domains based on intent classes, intent-slot combinations and example transcripts and parses. Experiments show that using synthetic parse data with JAT for existing domains can improve SSP performance on STOP by 1.4 % absolute EM. Using synthetic parse data with TTS for a new held-out domain improves EM on STOP for the held out domain by 2.6% absolute.

1 Introduction

Spoken Language Understanding (SLU) is essential for many real-world applications today including conversational agents and virtual assistants.

Spoken Semantic Parsing (SSP) is the SLU task that involves transforming a recording to a machine-comprehensible parse tree (Wang et al., 2023a). End-to-end models (Arora et al., 2023) operate directly on speech while cascade models (Futami et al., 2023) generate a semantic parse based on the speech transcript. Two-pass deliberation models (Le et al., 2022) combine the best of both worlds, by using first-pass transcripts and speech embeddings to perform spoken semantic parsing within a second pass. However, training such models with supervision requires matched triplets of speech, transcript, and semantic parse. Annotating these triplets is expensive, which limits the size of training data, and consequently model performance.

The need for matched data can be alleviated by developing methods that can use text-only unpaired data. Text data (transcript-semantic parse) is more easily obtained than speech – either from existing textual corpora or by prompting Large Language Models (LLMs), and training models with a small amount of paired speech-text data and a large amount of unpaired text is useful. It is non-trivial to incorporate text-only data into end-to-end models because model outputs cannot be obtained without speech inputs. Prior work has explored the use of text data for speech recognition (Wang et al., 2020a; Toshniwal et al., 2018; Hori et al., 2019). External language models trained on text can be used to interpolate token prediction probabilities (Meng et al., 2022), but require additional memory, making them unsuitable for on-device applications. Coordinated learning methods (Chen et al., 2022; Sainath et al., 2023) project speech and text to a shared embedding space for speech recognition, but such models require significant amounts of paired speech-text data to learn robust mappings. The final class of work generates speech representations for unpaired speech - Joint Audio Text (JAT) (Kim et al., 2022) uses mean speech embeddings from paired data to represent unpaired

043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083

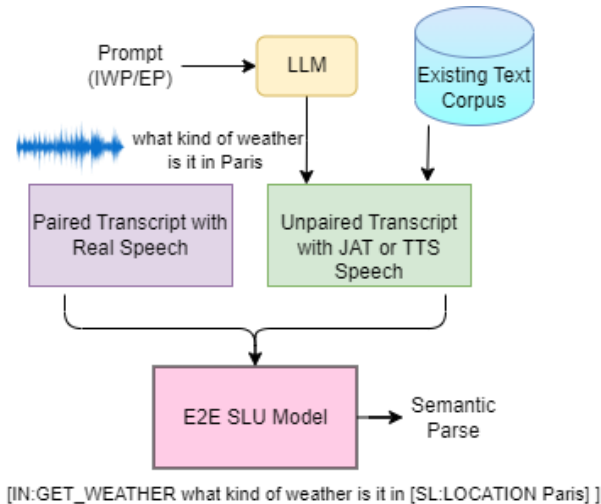


Figure 1: This paper: We describe ways to unpaired text to train deliberation models, where unpaired data can be obtained from LLMs or existing textual corpora. We use JAT or TTS to obtain speech representations of unpaired data

text. This is computationally inexpensive, but the speech embeddings do not contain information embedded in real speech. In contrast, synthetic speech from Text-to-speech (TTS) models (Wang et al., 2020a) produce informative speech representations, but they can be expensive to compute.

There are two cases where additional textual data may be acquired for semantic parsing – (a) to improve models on existing domains (ED) and (b) to support new domains (ND). In this paper, we compare JAT and TTS for SSP when unpaired text data is drawn from these two setups - ED and ND.

When unpaired text is not available from existing corpora, we propose to prompt Large Language Models (LLMs) (Ouyang et al., 2022; Touvron et al., 2023a,b) to generate textual data for SSP. LLMs have been used in prior work to generate synthetic data for text classification using approaches such as Self-Instruct (Wang et al., 2023b), Attr-Prompt (Yu et al., 2023a), ZeroGen (Ye et al., 2022), and more recently use in-context learning with seed samples (Yu et al., 2023b). Semantic parsing requires sequence labeling, i.e., (a) it requires the correct identification of identification of the number and identity of intent and slot tags, and (b) correct placement of entity and slot tags to form the right parse tree, all while not inserting unrelated or unseen intent and slot tags. Therefore, it is more complex to generate useful and diverse data for semantic parsing compared to other classification tasks.

Prior work (Tran and Tan, 2020) has proposed

the use of template-based masked training of BART to produce additional variants for masked words, however this limits the potential lexical diversity of the generated data, and requires significant amount of labeled data, which may not be available for the ND setting. Since LLMs can learn in-context and generalize better under few-shot settings, they consequently need fewer exemplars to generate diverse and high quality synthetic data for semantic parsing. In this paper, we address the task of generating synthetic text data for semantic parsing by using different prompting approaches with Llama 2.

For the ED setup, it is sufficient to generate transcripts (similar utterances) since semantic parses can be obtained from transcripts using pre-trained semantic parsers. We describe two prompting methods: (a) intent-word-based prompting (IWP), where the LLM produces transcripts corresponding to a particular intent class and containing words that co-occur with the intent, and (b) exemplar-based prompting (EP), where it generates transcripts that are similar to provided examples. We generate pseudo-labels for the generated utterances using a pre-trained RoBERTa (Liu et al., 2020) model and train SSP models using JAT. We find that EP is simpler but IWP generates the desired intent more often. Using data from both methods improves the Exact Match (EM) on STOP data by 1.4 points absolute.

For the ND setup, pre-trained models for pseudo-labeling are unavailable for the new domain(s), and hence LLMs are used to generate the seqlogical

242	4.1 Generating Textual Data for Existing	
243	Domains	
244	In the ED setup, we propose to use LLMs to gener-	
245	ate transcripts. Corresponding semantic parses are	
246	obtained using a pseudo-labeling textual semantic	
247	parse model trained on existing paired data. The	
248	semantic parse model here takes transcripts as in-	
249	puts and produces pseudo-label semantic parses as	
250	output. Transcripts can be generated using one of	
251	two prompting strategies, i.e., intent-word-based or	
252	exemplar-based.	
253	4.1.1 Intent Word-based prompting (IWP)	
254	The goal of IWP is to generate transcripts that may	
255	be classified under a certain intent, optionally con-	
256	taining "intent words". Intent words are the words	
257	from semantic parses that occur most frequently	
258	with given intents after removing stop-words. An	
259	example is shown in Figure 2. The 40 words that	
260	co-occur most frequently with every intent in the	
261	STOP data are used as intent words. 40 examples	
262	are generated for every intent and intent-word com-	
263	bination. Though IWP produces good synthetic	
264	data, it is limited by the fact that words that co-	
265	occur less frequently with the intent are less related	
266	to the intent. Such examples produced with less	
267	relevant intent words may not be classified under	
268	the desired intent class. This also limits the amount	
269	of synthetic data that can be generated since the	
270	LLM cannot generate many unique examples using	
271	a small number of intent-intent word combinations.	
272	4.1.2 Exemplar-based Prompting (EP)	
273	Since LLMs are strong in-context learners (Wei	
274	et al., 2022), an alternative approach is to prompt	
275	LLMs to generate transcripts based on examples.	
276	For every intent-slot combination, we provide up	
277	to 4 random example transcripts and ask the model	
278	to generate 60 more transcripts that are similar	
279	but have diverse sentence structures. An example	
280	prompt is shown in Fig 3. Though the resulting	
281	transcripts may not always correspond to the intent	
282	classes from which the examples are drawn, this	
283	method enables us to generate larger volumes of	
284	data without duplication.	
285	4.1.3 Semantic Parse generation and Quality	
286	Assessment	
287	Transcripts generated by LLMs are first normal-	
288	ized – written text is converted to spoken form,	
289	punctuation except apostrophes are removed and	
290	text is transformed into lower case. Semantic parse	
	pseudo-labels are obtained from these normalized	291
	transcripts using a strong RoBERTa-based seman-	292
	tic parser trained on STOP (EM=86.8). To assess	293
	data quality, we compare the intent in the obtained	294
	pseudo-labels to the intent in the prompt for IWP or	295
	the intent of the provided examples for EP. Intent	296
	Match Accuracy (IMA) is defined as the percent-	297
	age of times the intent of the pseudo-label matches	298
	the desired intent of the prompt.	299
	4.2 Generating Transcript-Semantic Parse for	300
	New Domains	301
	For new domains, paired data and pre-trained mod-	302
	els are not available, and therefore, we would need	303
	to directly generate pairs of transcript and semantic	304
	parse. One way to do this is to generate pairs of	305
	semantic parse and corresponding transcript using	306
	LLMs directly, however, maintaining consistency	307
	across generated parses and transcripts is challeng-	308
	ing for current LLMs. Another alternative is to	309
	generate only the seqlogical form of the semantic	310
	parse from the LLM and infer the transcript from	311
	the parse. The seqlogical form of the parse, unlike	312
	the decoupled form, comprises all the words in the	313
	transcript along with slot and intent tags. Therefore,	314
	the transcript can be obtained from the seqlogical	315
	parse merely by removing slot and intent tags.	316
	4.2.1 Exemplar-based Prompting	317
	We assume that (a) the intents and slots that must	318
	be recognized for the new domain are known, (b)	319
	the slots that may occur with every intent, i.e., the	320
	intent-slot combinations are known, and (c) some	321
	manually annotated examples for every intent-slot	322
	combination are known. Using this information,	323
	LLMs can be prompted as shown in Figure 4 to	324
	produce new seqlogical parses for a given intent-	325
	slot combinations. The prompt first describes the	326
	steps to generate a valid seqlogical parse and then	327
	presents up to 3 examples of seqlogical parses with	328
	the desired intent-slot combinations.	329
	4.2.2 Post-processing	330
	The generated seqlogical parses are checked for in-	331
	valid placement of brackets, and Out of Vocabulary	332
	(OOV) intents and slots. OOV intents were fixed	333
	by re-prompting the model to replace OOV intents	334
	with correct intents and replace any intents other	335
	than the first. Any OOV slots are removed while	336
	retaining corresponding slot words.	337

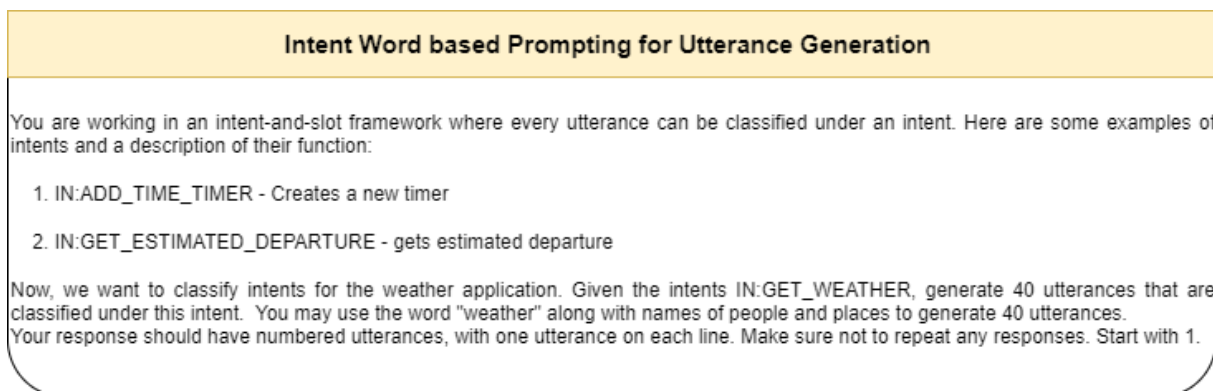


Figure 2: Example Prompt for IWP-based utterance generation

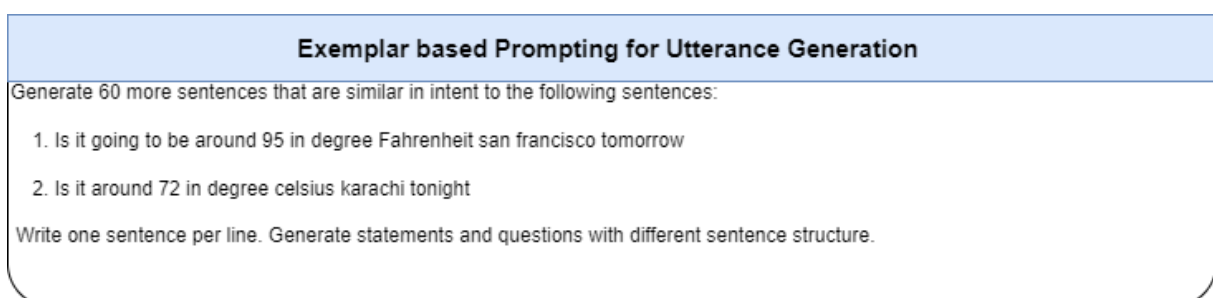


Figure 3: Example Prompt for EP-based utterance generation

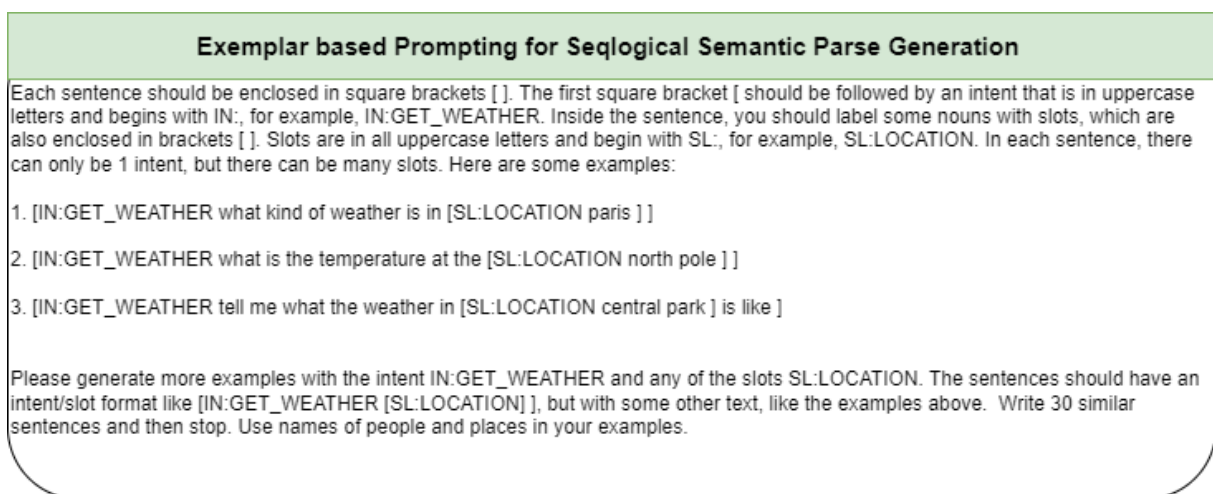


Figure 4: Example Prompt for EP-based generation of seqlogical parses

5 Experimental Setup

5.1 STOP Data, Model and Metrics

Data: STOP¹ (Tomasello et al., 2023) is a public dataset with 100 hours of real speech for spoken semantic parsing. STOP has data for 8 domains - alarm, event, messaging, music, navigation, reminder, timer, and weather. The data contains 28 unique intents and 82 slot types in all. Table 1 summarizes some statistics about the STOP dataset.

¹STOP was used in accordance with its LICENSE terms

Table 1: Dataset and Partition Statistics - STOP Dataset

Partition	Number of utterances
train	120,903
eval	33,380
test	75,617

Metrics: Exact Match (EM) is used to evaluate all our models. We report EM (No Err) and EM w/ Err, which are the Exact Match accuracies averaged over utterances with no ASR error and averaged

347
348
349
350

over utterances with any ASR error respectively.

Model Configuration: For the ASR module, we use RNNT with 3 layers of conformer in the encoder, 1 layer of LSTM in the predictor, and 1 linear layer in the joiner. For the deliberation model, we use attention in the Fusion module, 2 transformer encoder layers in the Pooling module, and a transformer decoder layer with a pointer-generator in the Decoder module (Kim et al., 2023). Models are optimized with Adam (Kingma and Ba, 2015), having a peak learning rate of $8e-3$.

Voicebox TTS Model: We use a Voicebox model trained on approximately 14k hours of manually transcribed data that comprises a diverse range of speakers, accents, topics, and acoustic conditions. The audio model has 12 transformer layers (Vaswani et al., 2017) containing 16 attention heads, convolutional positional embeddings (Baevski et al., 2020) and ALiBi self-attention bias (Press et al., 2021). Graphemes are embedded into 80-d features and concatenated with the 80-d log-mel features. The duration model has 8 transformer layers with 8 heads, and graphemes are embedded into 40-d features. Training hyperparameters are similar to the setup described in (Le et al., 2023).

Computational Cost : Our experiments were performed on a single node with 8 V100-32 GB GPUs on the cluster. Each run took approximately 18 hours for model training. For LLama2 inference, we used 4 x V100-32 GB or 2 x A100-40GB with model parallelism and fp32 precision. For Voicebox inference, we used 1X V100-32 GB GPUs over 40 parallel processes to speed up speech synthesis.

5.2 Setup: Textual Data from Text Corpora

For experiments where we assume textual data is available, we split the STOP datasets into two parts. We perform two experiments – one using the first and second splits as paired and unpaired data respectively and the other using the second and first splits as paired and unpaired data respectively. The average performance across these 2 experiments is reported in each case. In the ED setup, equal amounts of data from every domain are present in the two splits. For the ND setup, STOP is split by domain, where one split contains all training data from 4 domains (messaging, reminder, time, and weather), while the other split contains training data from the other 4 domains (alarm, event, music, and navigation). Both splits are designed to ensure that they have a nearly equal number of utterances.

5.3 Setup: Textual Data from LLMs

When unpaired data is not available, we use Llama 2.0 to generate examples for the ED and ND setups. For the ED setup, LLama 2.0 is used to generate utterances. We then use a pre-trained 12-layer RoBERTa model trained on STOP to generate pseudo-labels for the generated utterances. We augment STOP with the generated LLama 2.0 transcript-semantic parse. JAT is used to represent LLama 2 text.

For the ND setup, LLama 2.0 generated data is not suitable as a real test set since it does not have matching real speech. Therefore, we choose to partition the existing STOP data into 7 seen domains and 1 new domain - weather. We use exemplar-based prompting to generate transcript-semantic parse pairs for weather. For this, real examples of transcript-semantic parse from STOP are used. We use TTS to generate equivalent speech representations for the generated data. We compare the performance on the weather domain for models trained on (a) 7 domains of STOP, (b) 7 domains of STOP with examples for the weather (with TTS for examples and real speech for 7 domains), (c) 7 domains of STOP with examples and LLama 2.0 generated data, and (d) the topline that uses 7 domains of STOP with real data and TTS.

6 Experimental Results and Discussion

6.1 When textual data is available

Table 2 compares the performance of different models for the ED and ND settings where unpaired text is drawn from existing domains and new domains respectively. Across both ED and ND setups, we find that the use of unpaired text improves EM scores.

For the ED setup, we find that JAT and TTS achieve similar Exact Match scores. Since JAT is comparable in performance to TTS and relatively inexpensive compared to complex TTS models like Voicebox, JAT is optimal for the ED setup. TTS model training depends on the specific model, but in our case Voicebox training takes 3 days on 8 GPUs, and inference to produce synthetic speech takes 3 hours on 40 parallel GPU inference jobs. In comparison, JAT data preparation involves using mean speech embeddings, which takes 1 hour on 40 CPUs for the STOP training, evaluation and test data. Therefore, JAT indeed takes little time in comparison to TTS.

Further, the difference between JAT and TTS

Table 2: Comparing JAT and TTS as speech representations for unpaired text from ED and ND. Number of paired and unpaired utterances, and Exact Match (EM) is reported

	Model	#Pair/#Unpair	EM	EM(No Err)	EM w/ Err
ED	Baseline	60.4k / 0	64.25	80.51	24.37
	w/ JAT	60.4k / 60.4k	66.92	83.90	25.25
	w/ TTS	60.4 / 60.4k	67.05	83.88	25.80
ND	Baseline	60.7k / 0	33.28	41.32	13.54
	w/ JAT	60.7k / 60.1k	57.74	73.34	19.50
	w/ TTS	60.7k / 60.1k	63.95	80.70	22.88
	Topline	120.9k / 0	67.67	84.52	26.34

Table 3: Impact of Paired-Unpaired Data Ratio on JAT Performance under the Existing Domain Setting

Paired Data (%)	Unpaired Data (%)	EM-No Err	EM-ASR Error	EM (overall)
0	100	85.48	21.22	66.87
30	70	84.27	24.67	67.01
50	50	84.15	25.5	67.17
70	30	84.24	25.43	67.2
100	0	84.52	26.34	67.67

452 appears to be primarily on utterances with ASR er-
 453 rors, since synthetic speech representations can be
 454 used to reduce the impact of ASR errors on seman-
 455 tic parsing. For the ND setup, we find that though
 456 JAT outperforms the baseline, TTS outperforms
 457 JAT. This is because new domains may have dif-
 458 ferent entities and domain-specific terms that may
 459 be harder to recognize, and TTS provides valid
 460 speech representations that can be used to improve
 461 predictions based on the first-pass ASR. Figure 5
 462 shows that the amount of unpaired textual data is
 463 increased with constant paired data, relative gains
 464 increase to a point and saturate.

465 6.2 Ablation: Does JAT work for different 466 data ratios?

467 In this experiment, we vary the amount of paired
 468 data with speech-transcript-semantic parse and un-
 469 paired data with text only to analyze the impact on
 470 spoken semantic parsing performance.

471 From Table 3, we find that JAT works with only
 472 a 0.8 % degradation compared to the topline that
 473 uses 100% paired data in Exact Match even when
 474 no paired speech-text data is used. Therefore, this
 475 approach can generalize reasonably to other data ra-
 476 tios apart from the 50-50 ratio used in prior experi-
 477 ments. Further, utilizing more paired data improves
 478 performance on the cases when the transcript con-
 479 tains errors when compared to those where the
 480 transcript has no errors. This follows as a conse-
 481 quence of the fact that the transcript for unpaired
 482 text contains no ASR errors.

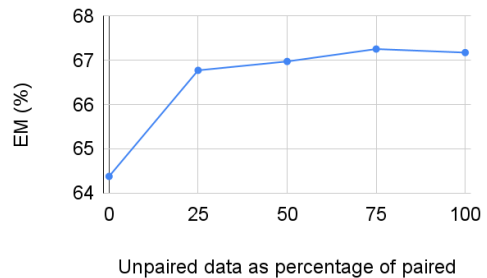


Figure 5: Impact of increasing unpaired text on EM

483 6.3 LLama 2.0 Generated Data: ED Setup

484 Table 4 compares various prompting strategies
 485 for generating utterances in the same domain us-
 486 ing Llama 2.0. We find that combining LLama-
 487 generated data with existing STOP data can im-
 488 prove performance across test examples with and
 489 without ASR errors. On further analysis, we find
 490 that significant improvements are observed across
 491 domains with relatively poor performance in the
 492 STOP baseline. Between IWP and EP, we find that
 493 EP is slightly better. Since EP is not constrained
 494 to generate utterances that may be classified under
 495 a given intent, the Intent Match Accuracy (IMA)
 496 is lower than that of IWP. Combining the data gen-
 497 erated from both these strategies further improves
 498 performance over the STOP baseline.

499 6.4 LLama 2.0 Generated Data: ND Setup

500 Table 5 compares the performance of baseline mod-
 501 els that have no data for weather or 360 examples
 502 for weather with models that use LLama 2.0 gen-

Table 4: Assessing the impact of augmenting the training data with LLama 2.0 generated utterances and RoBERTa pseudo-labels. EM is Exact Match Accuracy

Model	#Utts	IMA	EM	EM(No Err)	EM w/ Err
STOP Baseline	160k	-	67.37	84.52	26.34
+ IWP-JAT	230k	68.87	68.12	84.96	26.82
+ EP-JAT	218k	64.24	68.21	85.01	27.04
+ (IWP+EP)-JAT	298k	67.87	68.75	85.82	26.86

Table 5: Using TTS to generate speech for LLama 2.0 text when unpaired text is in an unseen new domain

Model	#Utts(Weather)	Weather EM	Overall EM
STOP 7 dom.		0	54.61
+ 3 real example-TTS		360	61.80
+ Exemplar LLama2-TTS		2,910	62.29
Topline: STOP Weather-TTS		2,910	66.33

erated data. Llama 2 generated text can improve performance by over 2 points absolute EM but lags behind the performance of a topline that uses data from STOP.

6.5 Challenges of using LLMs for generating large-scale data

While large language models can generate useful data based on the prompting strategies employed, there are certain challenges with generating large scale data, i.e., something of the order of few to many thousands of utterances.

LLMs can be reasonably consistent while responding within the current turn, but tend to repeat previously proposed examples after around 40 examples per input prompt, with variance arising from the complexity of different semantic parse structures. Due to input context limits while training, there is a limited number of unique and useful examples that can be elicited for every input prompt. It could be argued that each prompt can be presented multiple times with slight variations to obtain more data. However, LLMs are often not consistent across turns and end up repeating synthetic examples. One solution to this challenge could potentially involve using the "chat" formulation, where previous prompts and responses are part of the hidden states the model can attend to while producing new responses. However, due to memory limits, it is challenging to retain very long contexts in input memory, inhibiting the production of truly large scale data.

In this paper, we attempted to sample multiple times using different temperatures and seeds for every prompt to attempt to scale the obtained

data. This remains an interesting problem for future work.

7 Conclusion

We address the high cost of manually labeling speech-transcript-semantic parse data for spoken semantic parsing by enabling models to use text-only data. JAT is preferred for unpaired text in existing domains for its efficiency and gain of 2.5 % EM over a paired data baseline while remaining within 0.1 % EM of the more computationally expensive TTS. For unpaired text in new domains, TTS outperforms JAT by 6 % absolute EM overall, with a gain of 30.6 % EM over a paired baseline. When text data cannot be obtained from existing text corpora, we propose to prompt LLMs to generate transcript-semantic parse pairs. We show that using different prompting strategies, we can generate unpaired text data in relatively large volumes. Using JAT and TTS, we can leverage this LLM-generated data to further improve SSP by 1.4 % EM and 2.6 % EM absolute for existing and new domains.

559 Limitations

560 Our work uses the public open-source LLama2
561 LLM to generate synthetic data due to its open
562 source code, public model weights and determin-
563 istic generation behavior. However, prompting
564 behavior is not standard across all LLMs, and
565 though the general structure and strategy behind
566 our prompting can remain the same, specific and
567 small modifications may need to be made for dif-
568 ferent LLMs.

569 The STOP dataset, the only public dataset for
570 semantic parsing uses real but read speech, rather
571 than spontaneous speech. Making public data with
572 spontaneous speech and experimenting with such
573 will definitely be useful to explore.

574 Impact and Risks

575 Our work will enable the development of SLU mod-
576 els for tasks and languages where we have very
577 limited labelled data. We hope that this work also
578 spurs more collaboration across the fields of speech
579 and natural language processing, both of which are
580 needed to make progress in this area.

581 All the work in this paper was done in such a
582 manner so as to minimize the risk of misuse and
583 bias. Since the approach uses LLama to generate
584 synthetic data, potential risks include the perco-
585 lation of inherent biases in LLama into models
586 trained on such synthetic data.

587 References

588 Siddhant Arora, Hayato Futami, Shih-Lun Wu, Jes-
589 sica Huynh, Yifan Peng, Yosuke Kashiwagi, Emiru
590 Tsunoo, Brian Yan, and Shinji Watanabe. 2023. A
591 study on the integration of pipeline and e2e slu sys-
592 tems for spoken semantic parsing toward stop quality
593 challenge. In *ICASSP 2023-2023 IEEE International
594 Conference on Acoustics, Speech and Signal Process-
595 ing (ICASSP)*, pages 1–2. IEEE.

596 Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed,
597 and Michael Auli. 2020. wav2vec 2.0: A framework
598 for self-supervised learning of speech representations.
599 *Advances in neural information processing systems*,
600 33:12449–12460.

601 Zhehuai Chen, Yu Zhang, Andrew Rosenberg, Bhuvana
602 Ramabhadran, Pedro J. Moreno, Ankur Bapna, and
603 Heiga Zen. 2022. *MAESTRO: Matched Speech Text
604 Representations through Modality Matching*. In *Proc.
605 Interspeech 2022*, pages 4093–4097.

606 Hayato Futami, Jessica Huynh, Siddhant Arora, Shih-
607 Lun Wu, Yosuke Kashiwagi, Yifan Peng, Brian Yan,
608 Emiru Tsunoo, and Shinji Watanabe. 2023. The

609 pipeline system of asr and nlu with mlm-based data
610 augmentation toward stop low-resource challenge.
611 In *ICASSP 2023-2023 IEEE International Confer-
612 ence on Acoustics, Speech and Signal Processing
613 (ICASSP)*, pages 1–2. IEEE.

Alex Graves. 2012. Sequence transduction with
614 recurrent neural networks. *arXiv preprint
615 arXiv:1211.3711*. 616

Takaaki Hori, Ramon Astudillo, Tomoki Hayashi,
617 Yu Zhang, Shinji Watanabe, and Jonathan Le Roux.
618 2019. Cycle-consistency training for end-to-end
619 speech recognition. In *ICASSP 2019-2019 IEEE In-
620 ternational Conference on Acoustics, Speech and Sig-
621 nal Processing (ICASSP)*, pages 6271–6275. IEEE. 622

Suyoun Kim, Ke Li, Lucas Kabela, Rongqing Huang,
623 Jiedan Zhu, Ozlem Kalinli, and Duc Le. 2022. Joint
624 audio/text training for transformer rescorer of stream-
625 ing speech recognition. *EMNLP*. 626

Suyoun Kim, Yuan Shangguan, Jay Mahadeokar, An-
627 toine Bruguier, Christian Fuegen, Michael L Seltzer,
628 and Duc Le. 2021. Improved neural language model
629 fusion for streaming recurrent neural network trans-
630 ducer. In *ICASSP 2021-2021 IEEE International
631 Conference on Acoustics, Speech and Signal Process-
632 ing (ICASSP)*, pages 7333–7337. IEEE. 633

Suyoun Kim, Akshat Shrivastava, Duc Le, Ju Lin,
634 Ozlem Kalinli, and Michael L Seltzer. 2023. Modal-
635 ity confidence aware training for robust end-to-end
636 spoken language understanding. *Interspeech*. 637

Diederik P. Kingma and Jimmy Ba. 2015. *Adam: A
638 method for stochastic optimization*. In *3rd Inter-
639 national Conference on Learning Representations,
640 ICLR 2015, San Diego, CA, USA, May 7-9, 2015,
641 Conference Track Proceedings*. 642

Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020.
643 Hifi-gan: Generative adversarial networks for effi-
644 cient and high fidelity speech synthesis. *Advances in
645 Neural Information Processing Systems*, 33:17022–
646 17033. 647

Duc Le, Akshat Shrivastava, Paden Tomasello, Suy-
648 oun Kim, Aleksandr Livshits, Ozlem Kalinli, and
649 Michael L Seltzer. 2022. Deliberation model for on-
650 device spoken language understanding. *Interspeech*. 651

Matthew Le, Apoorv Vyas, Bowen Shi, Brian Kar-
652 rer, Leda Sari, Rashel Moritz, Mary Williamson,
653 Vimal Manohar, Yossi Adi, Jay Mahadeokar, et al.
654 2023. Voicebox: Text-guided multilingual uni-
655 versal speech generation at scale. *arXiv preprint
656 arXiv:2306.15687*. 657

Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu,
658 Maximilian Nickel, and Matt Le. 2022. Flow
659 matching for generative modeling. *arXiv preprint
660 arXiv:2210.02747*. 661

C. Liu, F. Zhang, D. Le, S. Kim, Y. Saraf, and G. Zweig.
662 2021. Improving RNN Transducer Based ASR with
663 Auxiliary Tasks. In *Proc. SLT*. 664

665	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	Ke Tran and Ming Tan. 2020. Generating synthetic data	721
666	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	for task-oriented semantic parsing with hierarchical	722
667	Luke Zettlemoyer, and Veselin Stoyanov. 2020.	representations . In <i>Proceedings of the Fourth Work-</i>	723
668	Ro{bert}a: A robustly optimized {bert} pretraining	<i>shop on Structured Prediction for NLP</i> , pages 17–21,	724
669	approach .	Online. Association for Computational Linguistics.	725
670	Zhong Meng, Yashesh Gaur, Naoyuki Kanda, Jinyu Li,	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	726
671	Xie Chen, Yu Wu, and Yifan Gong. 2022. Internal	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	727
672	Language Model Adaptation with Text-Only Data	Kaiser, and Illia Polosukhin. 2017. Attention is all	728
673	for End-to-End Speech Recognition . In <i>Proc. Inter-</i>	you need .	729
674	<i>speech 2022</i> , pages 2608–2612.	Gary Wang, Andrew Rosenberg, Zhehuai Chen,	730
675	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	Yu Zhang, Bhuvana Ramabhadran, Yonghui Wu, and	731
676	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	Pedro Moreno. 2020a. Improving speech recognition	732
677	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	using consistent predictions on synthesized speech .	733
678	2022. Training language models to follow instruc-	In <i>ICASSP 2020 - 2020 IEEE International Confer-</i>	734
679	tions with human feedback. <i>Advances in Neural</i>	<i>ence on Acoustics, Speech and Signal Processing</i>	735
680	<i>Information Processing Systems</i> , 35:27730–27744.	(<i>ICASSP</i>), pages 7029–7033.	736
681	Ofir Press, Noah A Smith, and Mike Lewis. 2021.	Peidong Wang, Tara N Sainath, and Ron J Weiss.	737
682	Train short, test long: Attention with linear biases	2020b. Multitask training with text data for	738
683	enables input length extrapolation. <i>arXiv preprint</i>	end-to-end speech recognition. <i>arXiv preprint</i>	739
684	<i>arXiv:2108.12409</i> .	<i>arXiv:2010.14318</i> .	740
685	Tara N Sainath, Ruoming Pang, Ron J Weiss, Yanzhang	Sid Wang, Akshat Shrivastava, and Sasha Livshits.	741
686	He, Chung-cheng Chiu, and Trevor Strohman. 2020.	2023a. Treepiece: Faster semantic parsing via tree	742
687	An attention-based joint acoustic and text on-device	tokenization .	743
688	end-to-end model. In <i>ICASSP 2020-2020 IEEE Inter-</i>	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa	744
689	<i>national Conference on Acoustics, Speech and Signal</i>	Liu, Noah A. Smith, Daniel Khoshabi, and Hannaneh	745
690	<i>Processing (ICASSP)</i> , pages 7039–7043. IEEE.	Hajishirzi. 2023b. Self-instruct: Aligning language	746
691	Tara N Sainath, Rohit Prabhavalkar, Ankur Bapna,	models with self-generated instructions . In <i>Proceed-</i>	747
692	Yu Zhang, Zhouyuan Huo, Zhehuai Chen, Bo Li,	<i>ings of the 61st Annual Meeting of the Association for</i>	748
693	Weiran Wang, and Trevor Strohman. 2023. Joist: A	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	749
694	joint speech and text streaming model for asr. In	pages 13484–13508, Toronto, Canada. Association	750
695	<i>2022 IEEE Spoken Language Technology Workshop</i>	for Computational Linguistics.	751
696	(<i>SLT</i>), pages 52–59. IEEE.	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,	752
697	Paden Tomasello, Akshat Shrivastava, Daniel Lazar, Po-	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,	753
698	Chun Hsu, Duc Le, Adithya Sagar, Ali Elkahky, Jade	Maarten Bosma, Denny Zhou, Donald Metzler, Ed H.	754
699	Copet, Wei-Ning Hsu, Yossi Adi, et al. 2023. Stop: A	Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy	755
700	dataset for spoken task oriented semantic parsing. In	Liang, Jeff Dean, and William Fedus. 2022. Emer-	756
701	<i>2022 IEEE Spoken Language Technology Workshop</i>	gent abilities of large language models . <i>Transactions</i>	757
702	(<i>SLT</i>), pages 991–998. IEEE.	<i>on Machine Learning Research</i> . Survey Certifica-	758
703	Shubham Toshniwal, Anjuli Kannan, Chung-Cheng	tion.	759
704	Chiu, Yonghui Wu, Tara N Sainath, and Karen	Jiacheng Ye, Jiahui Gao, Qintong Li, Hang Xu, Jiangtao	760
705	Livescu. 2018. A comparison of techniques for lan-	Feng, Zhiyong Wu, Tao Yu, and Lingpeng Kong.	761
706	guage model integration in encoder-decoder speech	2022. ZeroGen: Efficient zero-shot learning via	762
707	recognition. In <i>2018 IEEE spoken language technol-</i>	dataset generation . In <i>Proceedings of the 2022 Con-</i>	763
708	<i>ogy workshop (SLT)</i> , pages 369–375. IEEE.	<i>ference on Empirical Methods in Natural Language</i>	764
709	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	<i>Processing</i> , pages 11653–11669, Abu Dhabi, United	765
710	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	Arab Emirates. Association for Computational Lin-	766
711	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	guistics.	767
712	Azhar, et al. 2023a. Llama: Open and efficient	Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng,	768
713	foundation language models. <i>arXiv preprint</i>	Alexander Ratner, Ranjay Krishna, Jiaming Shen,	769
714	<i>arXiv:2302.13971</i> .	and Chao Zhang. 2023a. Large language model as	770
715	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	attributed training data generator: A tale of diversity	771
716	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	and bias .	772
717	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	Yue Yu, Yuchen Zhuang, Rongzhi Zhang, Yu Meng,	773
718	Bhosale, et al. 2023b. Llama 2: Open founda-	Jiaming Shen, and Chao Zhang. 2023b. ReGen:	774
719	tion and fine-tuned chat models. <i>arXiv preprint</i>	Zero-shot text classification via training data genera-	775
720	<i>arXiv:2307.09288</i> .	tion with progressive dense retrieval . In <i>Findings of</i>	776
		<i>the Association for Computational Linguistics: ACL</i>	777

778 2023, pages 11782–11805, Toronto, Canada. Associ-
779 ation for Computational Linguistics.

780 **A Example Appendix**

781 This is an appendix.