

Using large language models (LLMs) to synthesize training data

Prompt engineering enables researchers to generate customized training examples for lightweight “student” models.

The machine learning models that power conversational agents like Alexa are typically trained on labeled data, but data collection and labeling are expensive and complex, creating a bottleneck in the development process.

Large language models (LLMs) such as the 20-billion-parameter Alexa Teacher Model (AlexaTM 20B) might look like a way to break that bottleneck, since they excel in few-shot settings — i.e., when only a handful of labeled examples are available. But their size and computational costs are unsuitable for runtime systems, which require low latency and support high traffic volumes.

To enable models that are lightweight enough for runtime use, even when real training data is scarce, we propose Teaching via data (Tvd), in which we use an LLM-based “teacher” model to generate synthetic training data for a specific task, then use the generated data to fine-tune a smaller “student” model.

This blog post covers two of our recent papers on Tvd, LINGUIST, published at the 2022 International Conference on Computational Linguistics (COLING), generates training data for joint intent classification and slot tagging (IC+ST). CLASP, published at the 2022 Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (AAACL), generates training data for semantic parsing. Both tasks are core components of conversational AI.

We show that LINGUIST data generation improves on popular multilingual IC+ST benchmarks by 2 to 4 points absolute, while CLASP data generation improves multilingual semantic parsing by 5 to 6 points absolute.

The AlexaTM 20B model used in CLASP is now available on AWS JumpStart.

LINGUIST

Conversational-AI agents use intent classification and slot tagging (IC+ST) to understand the intent of a speaker’s request and identify the entities relevant to fulfilling that request. For example, when an agent is asked to “play ‘Wake Me Up’ by Avicii”, it might identify the intent as PlayMusic, with the slot value “wake me up” assigned to the slot Song and “Avicii” assigned to Artist. (Slot tagging in this context is also known as named-entity recognition, or NER.)

With real-world agents, the set of intents and slots grows over time as developers add support for new use cases. Furthermore, multilingual agents such as Alexa seek to maintain parity across languages when new intents and slots are developed, creating an additional bottleneck during development.

Suppose, for example, that we’re enabling a multilingual agent to understand the new intent GetWeather. To begin with, the intent may have only two associated utterances, in English and no other languages, annotated with the slots City and DayOfWeek. These two utterances alone are not enough to build a strong multilingual IC+ST model, so we need to obtain more training data.

Example: new GetWeather intent, two starter utterances:

- what’s the weather in [City: boston] on [DayOfWeek: monday]
- tell me the forecast for [DayOfWeek: friday] in [City: seattle] please

Sample starter utterances for the GetWeather intent.

A simple baseline approach to expanding this dataset to a new language is to translate the text. Here is an example using AlexaTM 20B with an in-context one-shot prompt. The text in the yellow box is the input to the model, and we can sample as many outputs from the model as we want, shown in the blue boxes.

One-Shot Translation Input

[CLM] Sentence: play the song wake me up by avicci

 Translation in French: joue la chanson wake me up par avicci
 Sentence: what’s the weather in boston on monday
 Translation in French:

AlexaTM 20B Model

Generated Output Examples

quel temps fera-t-il à boston le lundi
est-ce qu'il fait beau à boston lundi
comment est la météo à boston lundi

Alternate translations sampled from AlexaTM 20B.

To get more examples in the original English, we can either translate these French outputs back to English (back-translation) or directly use a paraphrasing model, such as, again, AlexaTM 20B with an in-context prompt:

One-Shot Paraphrase Input

[CLM] another way to say "play wake me up by avicci" is "can you play that avicci song called wake me up"
 another way to say "what's the weather in boston on monday" is "

AlexaTM 20B Model

Generated Output Examples

how’s the weather in boston on monday
what will the weather be like in boston on monday
will the weather be good on monday in boston

Using AlexaTM 20B as a paraphrase generator.

While these approaches go a long way, they have two key limitations: (1) the outputs don’t have the slot tags labeled, so we need to use a separate model (e.g., one that does word alignment) to guess which output words are City and which DayOfWeek, a process that introduces noise, and (2) we cannot control the outputs — say, by restricting them to specific slot types and values.

limitation

① 新的数据没有 slot tags
→ word alignment task
→ noise

② 不可控输出

对话历史是更新的，与 Alexa 一起使用。

输入语句 + 对话历史 → NER → 行动 API/NO → 行动参数填充 → 预测 Alexa 的回合 → 是 → 模型预测回合结束。

等待用户输入。

通过对话历史和输入语句，模型预测下一个动作。

如果模型预测回合结束，则返回给用户；否则，返回到步骤 1。

如果模型预测回合结束，则返回给