

ConveRT: Efficient and Accurate Conversational Representations from Transformers

github.com/PolyAI-LDN/polyai-models

Matthew Henderson, Iñigo Casanueva, Nikola Mrkšić,
Pei-Hao Su, Tsung-Hsien Wen and Ivan Vulić

matt@poly-ai.com

PolyAI Limited, London, UK

Abstract

General-purpose pretrained sentence encoders such as BERT are not ideal for real-world conversational AI applications; they are computationally heavy, slow, and expensive to train. We propose **ConveRT** (**C**onversational **R**epresentations from **T**ransformers), a pre-training framework for conversational tasks satisfying all the following requirements: it is effective, affordable, and quick to train. We pretrain using a retrieval-based response selection task, effectively leveraging quantization and subword-level parameterization in the dual encoder to build a lightweight memory- and energy-efficient model. We show that ConveRT achieves state-of-the-art performance across widely established response selection tasks. We also demonstrate that the use of extended dialog history as context yields further performance gains. Finally, we show that pretrained representations from the proposed encoder can be transferred to the intent classification task, yielding strong results across three diverse data sets. ConveRT trains substantially faster than standard sentence encoders or previous state-of-the-art dual encoders. With its reduced size and superior performance, we believe this model promises wider portability and scalability for Conversational AI applications.

1 Introduction

Dialog systems, also referred to as conversational systems or conversational agents, have found use in a wide range of applications. They assist users in accomplishing well-defined tasks such as finding and booking restaurants, hotels, and flights (Hemphill et al., 1990; Williams, 2012; El Asri et al., 2017), with further use in tourist information (Budzianowski et al., 2018), language learning (Raux et al., 2003; Chen et al., 2017), entertainment (Fraser et al., 2018), and healthcare (Laranjo et al., 2018; Fadhil and Schiavo, 2019). They are also key

components of intelligent virtual assistants such as Siri, Alexa, Cortana, and Google Assistant.

Data-driven task-oriented dialog systems require domain-specific labelled data: annotations for intents, explicit dialog states, and mentioned entities (Williams, 2014; Wen et al., 2017b,a; Ramadan et al., 2018; Liu et al., 2018; Zhao et al., 2019b). This makes the scaling and maintenance of such systems very challenging. Transfer learning on top of pretrained models (Devlin et al., 2019; Liu et al., 2019, *inter alia*) provides one avenue for reducing the amount of annotated data required to train models capable of generalization.

Pretrained models making use of language-model (LM) based learning objectives have become prevalent across the NLP research community. When it comes to dialog systems, *response selection* provides a more suitable pretraining task for learning representations that can encapsulate conversational cues. Such models can be pretrained using large corpora of natural unlabelled *conversational* data (Henderson et al., 2019b; Mehri et al., 2019). Response selection is also directly applicable to retrieval-based dialog systems, a popular and elegant approach to framing dialog (Wu et al., 2017; Weston et al., 2018; Mazaré et al., 2018; Gunasekara et al., 2019; Henderson et al., 2019b).¹

Response Selection is a task of selecting the most appropriate *response* given the dialog history (Wang et al., 2013; Al-Rfou et al., 2016; Yang et al., 2018; Du and Black, 2018; Chaudhuri et al., 2018). This task is central to retrieval-based dialog systems, which typically encode the *context* and a

¹Retrieval-based dialog is popular because posing dialog as response selection (Gunasekara et al., 2019) simplifies system design (Boussaha et al., 2019). Unlike modular or end-to-end task-oriented systems, retrieval-based ones do not rely on dedicated modules for language understanding, dialog management, and generation. They mitigate the requirements for explicit task-specific semantics hand-crafted by domain experts (Henderson et al., 2014; Mrkšić et al., 2015, 2017).

large collection of responses in a joint semantic space, and then retrieve the most relevant response by matching the query representation against the encodings of each candidate response. The key idea is to: **1)** make use of large unlabelled conversational datasets (such as Reddit conversational threads) to *pretrain* a neural model on the general-purpose response selection task; and then **2)** *fine-tune* this model, potentially with additional network layers, using much smaller amounts of task-specific data.

Dual-encoder architectures pretrained on response selection have become increasingly popular in the dialog community (Cer et al., 2018; Humeau et al., 2020; Henderson et al., 2019b). In recent work, Henderson et al. (2019a) show that standard pretraining LM-based architectures cannot match the performance of dual encoders when applied to dialog tasks such as response retrieval.

Scalability and Portability. A fundamental problem with pretrained models is their large number of parameters (see Table 2 later): they are typically highly computationally expensive to both train and run (Liu et al., 2019). Such high memory footprints and computational requirements hinder quick deployment as well as their wide portability, scalability, and research-oriented exploration. The need to make pretrained models more compact has been recognized recently, with a line of work focused on building more efficient pretraining and fine-tuning protocols (Tang et al., 2019; Sanh et al., 2019). The desired reductions have been achieved through techniques such as distillation (Sanh et al., 2019), quantization-aware training (Zafir et al., 2019), weight pruning (Michel et al., 2019) or weight tying (Lan et al., 2019). However, the primary focus so far has been on optimizing the LM-based models, such as BERT.

ConveRT. This work introduces a *more compact pretrained response selection model* for dialog. ConveRT is only 59MB in size, making it significantly smaller than the previous state-of-the-art dual encoder (444MB). It is also more compact than other popular sentence encoders, as illustrated in Table 2. This notable reduction in size and training acceleration are achieved through combining 8-bit embedding quantization and quantization-aware training, subword-level parameterization, and pruned self-attention. Furthermore, the lightweight design allows us to reserve additional parameters to improve the expressiveness of the dual-encoder architecture; this leads

to *improved learning of conversational representations* that can be transferred to other dialog tasks (Casanueva et al., 2020; Bunk et al., 2020).

Multi-Context Modeling. ConveRT moves beyond the limiting single-context assumption made by Henderson et al. (2019b), where only the immediate preceding context was used to look for a relevant response. We propose a multi-context dual-encoder model which combines the immediate context with previous dialog history in the response selection task. The multi-context ConveRT variant remains compact (73MB in total), while offering improved performance on a range of established response selection tasks. We report significant gains over the previous state-of-the-art on benchmarks such as Ubuntu DSTC7 (Gunasekara et al., 2019), AmazonQA (Wan and McAuley, 2016) and Reddit response selection (Henderson et al., 2019a), both in single-context and multi-context scenarios. Moreover, we show that sentence encodings learned by the model can be transferred to other dialog tasks, reaching strong intent classification performance over three evaluation sets. Pretrained dual-encoder models, both single-context and multi-context ones, are shared as TensorFlow Hub modules at github.com/PolyAI-LDN/polyai-models.²

2 Methodology

Pretraining on Reddit Data. We assume working with English throughout the paper. Simplifying the conversational learning task to response selection, we can relate target dialog tasks to general-domain conversational data such as Reddit (Al-Rfou et al., 2016). This allows us to fine-tune the parameters of the task-specific response selection model, starting from the general-domain response selection model pretrained on Reddit. Similar to Henderson et al. (2019b), we choose Reddit for pretraining due to: **1)** its organic conversational structure; and **2)** its unmatched size, as the public repository of Reddit data comprises 727M (*input, response*) pairs.³

²Finally, our more compact neural response selection architecture is well aligned with the recent socially-aware initiatives on reducing costs and improving fairness and inclusion in NLP research and practice (Strubell et al., 2019; Mirzadeh et al., 2019; Schwartz et al., 2019). Cheaper training (pretraining the proposed dual-encoder model on the entire Reddit costs only 85 USD) and quicker development cycles offer new opportunities for more researchers and practitioners to tap into the construction of neural task-based dialog systems.

³github.com/PolyAI-LDN/conversational-datasets

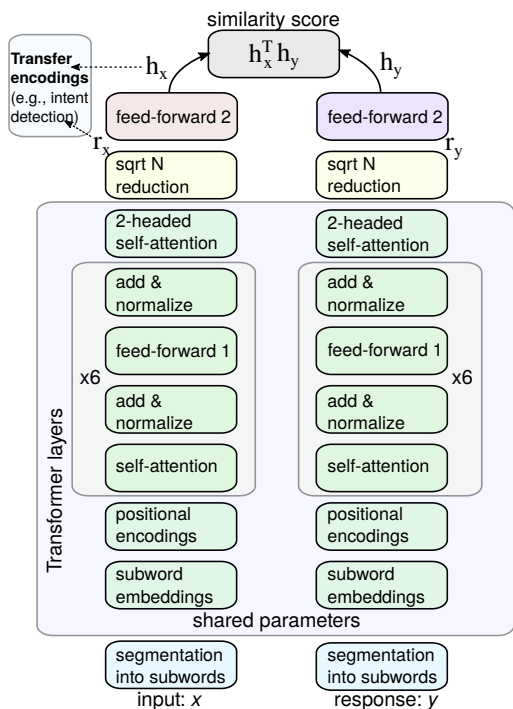


Figure 1: *Single-context ConveRT* dual-encoder model architecture. Its multi-context extension is illustrated in Figure 2. It is possible to *transfer* learned encodings at different network layers (e.g., r_x or the final h_x) to other tasks such as intent detection or value extraction (see §4). Note that the model uses two different feed-forward network (FFN) layers: 1) *feed-forward 1* is the standard FFN layer also used by Vaswani et al. (2017), and 2) *feed-forward 2* contains 3 fully-connected non-linear feed-forward layers followed by a linear layer which maps to the final encodings h_x and h_y (note that the two *feed-forward 2* networks do not share parameters, while the *feed-forward 1* parameters are shared).

2.1 More Compact Response Selection Model

We propose **ConveRT** – Conversational Representations from Transformers – a *compact dual-encoder pretraining architecture*, leveraging subword representations, transformer-style blocks, and quantization, as illustrated in Figure 1. ConveRT satisfies all the following requirements: it is effective, affordable, and quick to train.

Input and Response Representation. Prior to training, we obtain a vocabulary of subwords V shared by the input side and the response side: we randomly sample and lowercase 10M sentences from Reddit, and then iteratively run any subword tokenization algorithm.⁴ The final vocabulary V

⁴In the actual implementation, we use the same subword tokenization as Vaswani et al. (2018). We run it for 4 iterations and retain only subwords occurring at least 250 times, containing no more than 20 UTF8 characters, also disallowing more than 4 consecutive digits.

contains 31,476 subword tokens. During training and inference, if we encounter an OOV character it is treated as a subword token, where its ID is computed using a hash function, and it gets assigned to one of 1,000 additional “buckets” reserved for the OOVs. We therefore reserve parameters (i.e., embeddings) for the 31,476 subwords from V and for the additional 1,000 OOV-related buckets. At training and inference, after the initial word-level tokenization on UTF8 punctuation and word boundaries, input text x is split into subwords following a simple left-to-right greedy prefix matching (Vaswani et al., 2018). We tokenize all responses y during training in exactly the same manner.

Input and Response Encoder Networks. The subword embeddings then go through a series of transformations on both the input and the response side. The transformations are based on the standard Transformer architecture (Vaswani et al., 2017). Before going through the self-attention blocks, we add positional encodings to the subword embedding inputs. Previous work (e.g., BERT and related models) (Devlin et al., 2019; Lan et al., 2019, *inter alia*) learns a fixed number of positional encodings, one for each position in the sequence, allowing the model to represent a fixed number of positions. Instead, we learn two positional encoding matrices of different sizes- M^1 of dimensionality [47, 512] and M^2 of dimensionality [11, 512]. An embedding at position i is added to: $M_i^1 \bmod 47 + M_i^2 \bmod 11$.⁵

The next layers closely follow the original Transformer architecture with some notable differences. First, we set maximum relative attention (Shaw et al., 2018) in the six layers to the following respective values: [3, 5, 48, 48, 48, 48].⁶ This also helps the architecture to generalize to long sequences and distant dependencies: earlier layers are forced to group together meanings at the phrase level before later layers model larger patterns. We use single-headed attention throughout the network.⁷

Before going into a softmax, we add a bias to the attention scores that depends only on the rel-

⁵Note that since 47 and 11 are coprime, this gives $47 \cdot 11 = 517$ different possible positional encodings. Similar to the original (non-learned) positional encodings from Vaswani et al. (2017), the rationale behind this choice of positional encoding is to allow the model to generalize to unseen sequence lengths.

⁶We zero out in training and inference the attention scores for pairs of words if they are further apart than the set maximum relative attention values.

⁷Multi-headed attention requires running computations on 4-tensors: [batch, time, head, embedding], while for single-headed attention, this reduces to 3-tensors, and effectively speeds up training without hurting performance.

ative positions: $\alpha_{ij} \rightarrow \alpha_{ij} + B_{n-i+j}$ where B is a learned bias vector. This helps the model understand relative positions, but is much more computationally efficient than computing full relative positional encodings (Shaw et al., 2018). Again, it also helps the model generalize to longer sequences.

Six Transformer blocks use a 64-dim projection for computing attention weights, a 2,048-dim kernel (*feed-forward 1* in Figure 1), and 512-dim embeddings. Note that all Transformer layers use parameters that are fully shared between the input side and the response side. As in the Universal Sentence Encoder (USE) (Cer et al., 2018), we use square-root-of-N reduction to convert the embedding sequences to fixed-dimensional vectors. Two self-attention heads each compute weights for a weighted sum, which is scaled by the square root of the sequence length; the length is computed as the number of constituent subwords.⁸ The outputs of the reduction layer, labelled r_x and r_y in Figure 1, are 1,024-dimensional vectors that are fed to the two “side-specific” (i.e., they do not share parameters) feed-forward networks.

In other words, the vectors r_x and r_y go through a series of N_f l -dim feed-forward hidden layers ($N_f = 3$; $l = 1,024$) with skip connections, layer normalization, and orthogonal initialization. The activation function used in these networks and throughout the architecture is the fast GeLU approximation (Hendrycks and Gimpel, 2016): $GeLU(x) = x\sigma(1.702x)$. The final layer is linear and maps the text into the final L2-normalized 512-dim representation: h_x for the input text, and h_y for the corresponding response text (Figure 1).

Input-Response Interaction. The relevance of each response to the given input is then quantified by the score $S(x, y)$, computed as cosine similarity with annealing between the encodings h_x and h_y . It starts at 1 and ends at \sqrt{d} , linearly increasing over the first 10K training batches. Training proceeds in batches of K (*input, response*) pairs $(x_1, y_1), \dots, (x_K, y_K)$. The aim of the objective is to distinguish between the true relevant response (y_i) and irrelevant responses (i.e., negative samples) $y_j, j \neq i$ for each input sentence x_i . The training objective for a single batch of K pairs is as follows:

⁸In fact, rather than computing the self-attended sequence, then reducing it, we reduce the attention weights accordingly, and then directly apply them via matrix multiplication to the input sequence to get the final reduced representation, that is, we *fuse* these two operations. This is more computationally efficient, avoiding another 3-tensor multiplication.

$$J = \sum_{i=1}^K S(x_i, y_i) - \sum_{i=1}^K \log \sum_{j=1}^K e^{S(x_i, y_j)}.$$

The goal is to maximize the score of positive training pairs (x_i, y_i) and minimize the score of pairing each input x_i with K' negative examples, which are responses that are not associated with the input x_i : for simplicity, all other $K - 1$ from the current batch are used as negative examples.

Quantization. Very recent work has shown that large models of language can be made more compact by applying quantization techniques (Han et al., 2016): e.g., quantized versions of Transformer-based machine translation systems (Bhandare et al., 2019) and BERT (Shen et al., 2019; Zhao et al., 2019a; Zafrir et al., 2019) are now available. In this work, we focus on enabling quantization-aware conversational pretraining on the response selection task. We show that the dual-encoder ConveRT model from Figure 1 can be also be trained in a quantization-aware manner. Rather than the standard 32-bits per parameter, all embedding parameters are represented using only 8 bits, and other network parameters with just 16 bits; they are trained in a quantization-aware manner by adapting the mixed precision training scheme from Micikevicius et al. (2018). It keeps shadow copies of each variable with 32bit Floating Point (FP32) precision, but uses FP16-cast versions in the computations and inference models. Some operations in the graph, however, require FP32 precision to be numerically stable: layer normalization, L2-normalization, and softmax in attention layers.

Again, following Micikevicius et al. (2018), the final loss is scaled by 128, and the updates to the shadow FP32 variables are scaled back by 1/128: this allows the gradient computations to stay well represented by FP16 (e.g., they will not get rounded to zero). The subword embeddings are stored using 8-bits per parameter, and the quantization range is adjusted dynamically through training. It is updated periodically to contain all of the embedding values that have so-far been learned, with room for growth above and below - 10% of the range, or 0.01 - whichever is larger. Finally, quantization also allows doubling the batch size, which also has a favourable effect of increasing the number of negative examples in training.

Multi-Context ConveRT. Figure 1 depicts a single-context dual encoder architecture. Intuitively, the single-context assumption is limiting for modeling multi-turn conversations, where strong conversational cues can be found in earlier dialog

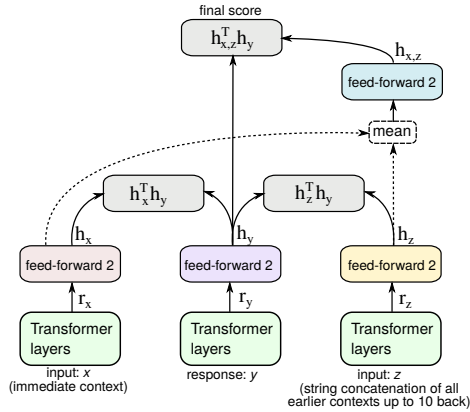


Figure 2: *Multi-context ConVeRT*. It models 1) the interaction between the immediate context and its accompanying response, 2) the interaction of the response with up to 10 earlier contexts from the conversation history, as well as 3) the interaction of the full context with the response. *Transformer layers* refer to the standard Transformer architecture also used in the single-context encoder model in Figure 1; the *feed-forward 2* blocks are the same as with the single-context encoder architecture, see Figure 1. The block *mean* refers to simple averaging of two context encodings h_x and h_z .

history, and there has been a body of work on leveraging richer dialog history for response selection (Chaudhuri et al., 2018; Zhou et al., 2018; Humeau et al., 2020). Taking a simple illustrative example:

Student: I’m very interested in representation learning.

Teacher: Do you have any experience in PyTorch?

Student: Not really.

Teacher: And what about TensorFlow?

Selecting the last Teacher’s response would be very difficult given only the immediate preceding context. However, the task becomes easier when taking into account the entire context of the conversation. We thus construct a *multi-context dual-encoder model* by using up to 10 more previous messages in a Reddit thread. The extra 10 contexts are concatenated from most recent to oldest, and treated as an extra feature in the network, as shown in Figure 2. Note that all context representations are still independent from the representation of a candidate response, so we can still do efficient response retrieval and training. The full training objective is a linear combination of three sub-objectives: 1) ranking responses given the immediate context (i.e., this is equal to the single-context model from §2.1), 2) ranking responses given only the extra (non-immediate) contexts, and 3) ranking responses given the averaged representation of the

immediate context and additional contexts.⁹

3 Experimental Setup

Training Data and Setup. We base all our (pre)training on the large Reddit conversational corpus (Henderson et al., 2019a) derived from 3.7B Reddit comments: it comprises 727M (*input, response*) pairs for single-context modeling – 654M pairs are reserved for training, the rest is used for testing. We truncate sequences to 60 subwords, embedding size is set to 512 for all subword embeddings and bucket embeddings, and the final encodings h_x , h_y , h_z , and $h_{x,z}$ are all 512-dimensional. The hidden layer size of *feed forward 2* networks is set to 1,024 (with $N_f = 3$ hidden layers used).

We train using ADADELTA with $\rho = 0.9$ (Zeiler, 2012), batch size of 512, and a learning rate of 1.0 annealed to 0.001 with cosine decay over training. L2-regularization of 10^{-5} is used, subword embedding gradients are clipped to 1.0, and label smoothing of 0.2 is applied.¹⁰

We pretrain the model on Reddit on 12 GPU nodes with one Tesla K80 each for 18 hours; this is typically sufficient to reach convergence. The total pretraining cost is roughly \$85 on Google Cloud Platform. This pretraining regime is orders of magnitude cheaper and more efficient than the prevalent pretrained NLP models such as BERT, GPT-2, XLNet, and RoBERTa (Strubell et al., 2019).

Baselines. We report results on the response selection tasks and compare against the standard set of baselines (Henderson et al., 2019a). First, we compare to a simple keyword matching baseline based on TF-IDF query-response scoring (Manning et al., 2008), and then with a representative sample of publicly available neural encoders that embed inputs and responses into a vector space relying on various pretraining objectives: (1) The larger

⁹Combining multiple objectives in a dual-encoder framework has also been done by Al-Rfou et al. (2016) and Henderson et al. (2017). Note that more sophisticated solutions to fusing dialog history are possible such as using attention over older contexts as done by Vlasov et al. (2019) on the much smaller MultiWOZ 2.1 dataset (Eric et al., 2019), but we have opted for simple concatenation as an efficient solution for training on the large Reddit data. The multiple objectives result in quicker learning, and also give useful diagnostic probes into the performance of each feature throughout training.

¹⁰The label smoothing technique (Szegedy et al., 2016) reduces overfitting by preventing a network to assign full probability to the correct training example (Pereyra et al., 2017). It means that each positive example in each batch is assigned the probability of 0.8, while the remaining probability mass is evenly redistributed across in-batch negative examples.

variant of Universal Sentence Encoder (Cer et al., 2018) (USE-LARGE); (2) The large variant of BERT (Devlin et al., 2019) (BERT-LARGE). We also compare to two recent dual-encoder architectures: (3) USE-QA is a dual question-answer encoder version of the USE (large) model (Chidambaram et al., 2019).¹¹ (4) POLYAI-DUAL is the best-performing dual-encoder model from Henderson et al. (2019b) pretrained on Reddit response selection. For baseline models 1-3, we report the results with the MAP response selection variant (Henderson et al., 2019a): it showed much stronger performance than a simpler similarity-based variant which directly ranks responses according to their cosine similarity with the context vector. MAP learns to (linearly) map the response vectors to the input vector space.

Response Selection: Evaluation Tasks. We report response selection performance on Reddit test set (Henderson et al., 2019a) with both single-context and multi-context ConveRT variants. For multi-context ConveRT, the averaged representation of (immediate and previous) context is used in evaluation. The models are applied directly on the Reddit test data without any further fine-tuning. We also evaluate on two other well-known response selection problems in different domains. (1) AMAZONQA (Wan and McAuley, 2016) is an e-commerce data set which contains information about Amazon products in the form of question-answer pairs: out of 3.6M (single-context) QA pairs, 300K pairs are reserved for testing. (2) DSTC7-UBUNTU is based on the Ubuntu v2 corpus (Lowe et al., 2017): it contains 1M+ conversations in a highly technical domain (i.e., Ubuntu technical support). DSTC7-UBUNTU uses 100K conversations for training, 10K for validation, and 5K conversations are used for testing (Gunasekara et al., 2019).

For DSTC7-UBUNTU we fine-tune for 60K training steps: it takes around 2h on 12 GPU workers. The learning rate starts at 0.1, and is annealed to 0.0001 using cosine decay over training. We use a batch size of 256, and dropout of 0.2 after the embedding and self-attention layers. We use the same fine-tuning regime for AMAZONQA. For DSTC7-UBUNTU, extra contexts are prepended with numerical strings 0–9 to help the model identify their position. We also release the fine-tuned models.

We evaluate with a standard IR-inspired eval-

¹¹Note that USE-QA encodes inputs/contexts and responses using separate sub-networks, while ConveRT (Figure 1) relies on full parameter sharing in the Transformer layers.

	# intents	# examples
Banking (customer service)	77	14.6K
Shopping (online shopping)	10	13.8K
Company FAQ	110	3.3K

Table 1: Intent classification data sets.

uation measure: $Recall@k$, used in prior work on retrieval-based dialog (Chaudhuri et al., 2018; Henderson et al., 2019b; Gunasekara et al., 2019). Given a set of N responses to the given input, where only one response is relevant, it indicates whether the relevant response occurs in the top k ranked candidates. We denote this measure as $\mathbf{R}_N@k$, and set $N = 100; k = 1: \mathbf{R}_{100}@1$.

Intent Classification: Task, Data, Setup. Pre-trained sentence encoders have become particularly popular due to the success of training models for downstream tasks on top of their learned representations, greatly improving the results compared to training from scratch, especially in low-data regimes (see Table 1). Therefore, we also probe the usefulness of ConveRT encodings for transfer learning in the intent classification task: the model must classify the user’s utterance into one of several predefined classes, that is, *intents* (e.g., within e-banking intents can be *card lost* or *replace card*). We use three internal intent classification datasets from three diverse domains, see Table 1, divided into train, dev and test sets using a 80/10/10 split.

We use the pretrained ConveRT encodings r_x on the input side (see Figure 1) as input to an intent classification model. We also experimented with later h_x encodings on the input side, but stronger results were observed with r_x . We train a 2-layer feed-forward net with dropout on top of r_x . SGD with a batch size of 32 is used, with early stopping after 5 epochs without improvement on the validation set. Layer sizes, dropout rate and learning rate are selected through grid search. We compare against two other standard sentence encoders again: USE-LARGE and BERT-LARGE. For ConveRT and USE-LARGE we keep the encoders fixed and train the classifier layers on top of the sentence encodings. For BERT-LARGE, we train on top of the CLS token and we fine-tune all its parameters.

4 Results and Discussion

Model Size, Training Time, Cost. Table 2 lists encoders from prior work along with their model size, and estimated model size after quantization. The reported numbers indicate the gains achieved

	Embedding parameters	Network parameters	Total size	Size after quantization
USE (Cer et al., 2018)	256 M	2 M	1033 MB	261 MB *
BERT-BASE (Devlin et al., 2019)	23 M	86 M	438 MB	196 MB */ 110 MB **
BERT-LARGE (Devlin et al., 2019)	31 M	304 M	1341 MB	639 MB */ 336 MB **
GPT-2 (Radford et al., 2019)	80 M	1462 M	6168 MB	3004 MB *
POLYAI-DUAL (Henderson et al., 2019b)	104 M	7 M	444 MB	118 MB
ConveRT (this work)	16 M	13 M	116 MB	59 MB

Table 2: Comparison of the proposed compact dual-encoder architecture for response selection to existing public standard sentence embedding models. (*) The size after quantization assumes embeddings can be quantized to 8 bits and network parameters to 16 bits, which has not been verified for the public models. (**) Best-case model size estimates of the BERT model after full 8-bit quantization based on the work of Zafir et al. (2019).

	Reddit	AmazonQA
TF-IDF	26.4	51.8
USE-LARGE-MAP	47.7	61.9
BERT-LARGE-MAP	24.0	44.1
USE-QA-MAP	46.6	70.7
POLYAI-DUAL	61.3	71.3
ConveRT (single-context)	68.2	84.3
ConveRT (multi-context)	71.8	–

Table 3: $\mathbf{R}_{100}@1 \times 100\%$ scores on Reddit test set and AMAZONQA. POLYAI-DUAL and ConveRT networks are fine-tuned on the training portion of AMAZONQA. Note that AMAZONQA by design supports only single-context response selection.

Model Configuration	
ConveRT	68.2
A: Multi-headed attention (8 64-dim heads)	68.5
B: No relative position bias	67.8
C: Without gradually increasing max attention span	67.7
D: Only 1 OOV bucket	68.0
E: 1-headed (instead of 2-headed) reduction	67.7
F: No skip connections in <i>feed forward 2</i>	67.8
D + E + F	66.7
B + C + D + E + F	66.6

Table 4: An ablation study illustrating the importance of different components in ConveRT: single-context response selection on Reddit ($\mathbf{R}_{100}@1$). Each experiment has been run for 966K steps (batch size 512).

through subword-level parameterization and quantization of ConveRT. Besides reduced training costs, ConveRT offers a reduced memory footprint and quicker training. We pretrain all our models for 18 hours only (on 12 16GB T4 GPUs), while a model compression technique DistilBERT (Sanh et al., 2019) (i.e., it reports $\approx 40\%$ relative reduction of the original BERT) trains on 8 16GB V100 GPUs for 90 hours, and larger models like RoBERTa require 1 full day of training on 1024 32GB V100 GPUs. The achieved size reduction and quick training also allow for quicker development and insightful ablation studies (see later in Table 4), and using quantization also improves training efficiency in terms of examples per second.

	$\mathbf{R}_{100}@1$	MRR
Best DSTC7 System	64.5	73.5
GPT*	48.9	59.5
BERT*	53.0	63.2
Bi-encoder (Humeau et al., 2020)	70.9	78.1
ConveRT (single-context)	38.2	49.2
ConveRT (multi-context)	71.2	78.8

Table 5: Results on DSTC7-UBUNTU. (*) Scores for GPT and BERT taken from Vig and Ramea (2019).

Response Selection on Reddit. The results are summarized in Table 3. Even single-context ConveRT achieves peak performance in the task, with substantial gains over the previous best reported score of Henderson et al. (2019b). It also substantially outperforms all the other models which were not pretrained directly on the response selection task, but on a standard LM task instead. The strongest baselines, however, are two dual-encoder architectures (i.e., USE-LARGE, USE-QA and POLYAI-DUAL); this illustrates the importance of explicitly distinguishing between inputs/contexts and responses when modeling response selection.

Table 3 also shows the importance of leveraging additional contexts (see Figure 2). Multi-context ConveRT achieves a state-of-the-art Reddit response selection score of **71.8%**. We observe similar benefits in other reported response selection tasks. We also note the results of 1) using only the sub-network that models the interaction between the immediate context and the response (i.e., the $h_x^T h_y$ interaction), and 2) artificially replacing the concatenated extra contexts z with an empty string. The respective scores are 65.7% and 65.6%. This suggests that multi-context ConveRT is also applicable to single-context scenarios when no extra contexts are provided for the target task.

Ablation Study. The efficient training regime also allows us to perform a variety of diagnostic experiments and ablations. We report results with variants of single-context ConveRT in Table 4. They

indicate that replacing single-headed with multi-headed attention leads to slight improvements, but this comes at a cost of slower (and consequently - more expensive) training. Using 1 instead of 1,000 OOV buckets leads only to a modest decrease in performance. Most importantly, the ablation study indicates that the final performance actually comes from the synergistic effect of applying a variety of components and technical design choices such as skip connections, 2-headed reductions, relative position biases, etc. While removing only one component at a time yields only modest performance losses, the results show that the loss adds up as we remove more components, and different components indeed contribute to the final score.¹²

Other Response Selection Tasks. The results on the AMAZONQA task are provided in Table 3. We see similar trends as with Reddit evaluation. Fine-tuned ConveRT reaches a new state-of-the-art score, and the strongest baselines are again dual-encoder networks. Fine-tuned POLYAI-DUAL, which was pretrained on exactly the same data, cannot match ConveRT’s performance.¹³

The results on DSTC7-UBUNTU are summarized in Table 5. First, they suggest very competitive performance of multi-context ConveRT model: it outperforms the best-scoring system from the official DSTC7 challenge (Gunasekara et al., 2019). It is an encouraging finding, given that multi-context ConveRT relies on simple context concatenation without any additional attention mechanisms. We leave the investigation of such more sophisticated models to integrate additional contexts for future work. Multi-context ConveRT can also match or even surpass the performance of another dual-encoder architecture from Humeau et al. (2020). Their dual encoder (i.e., *bi-encoder*) is based on the BERT-base architecture (Humeau et al., 2020): it relies on 12 Transformer blocks, 12 attention heads, and a hidden size dimensionality of 768 (while we use 512). Training with that model is roughly 5× slower, and

¹²Furthermore, quick development and short training times also allow us to treat some of the component choices as hyperparameter choices. It effectively means that such configuration choices can also be fine-tuned similar to any other hyperparameter to optimize the final retrieval performance.

¹³Interestingly, directly applying ConveRT to AMAZONQA without any fine-tuning also yields a reasonably high score of 67.0%. Moreover, learning the mapping function between inputs and responses (again without any fine-tuning) for ConveRT the same way as is done for USE-QA-MAP results in the score of 71.6%, which outperforms USE-QA-MAP (70.7%). The gap to the fine-tuned model’s performance, however, indicates the importance of in-domain fine-tuning.

	Banking	Shopping	Company FAQ
USE-LARGE	92.2	94.0	62.4
BERT-LARGE	93.2	94.3	61.2
ConveRT	92.7	94.5	64.3

Table 6: Intent classification results.

the pretraining objective is more complex: they use the standard BERT pretraining objective plus next utterance classification. Moreover, their model is trained on 32 v100 GPUs for 14 days, which makes it roughly 50× more expensive than ConveRT.

Intent Classification. The results are summarized in Table 6: we report the results of two strongest baselines for brevity. The scores show very competitive performance of ConveRT encodings r_x transferred to another dialog task. They outperform USE-LARGE in all three tasks and BERT-LARGE in 2/3 tasks. Note that, besides quicker pretraining, intent classifiers based on ConveRT encodings train 40 times faster than BERT-LARGE-based ones, as only the classification layers are trained for ConveRT. In sum, these preliminary results suggest that ConveRT as a sentence encoder can be useful beyond the core response selection task. The usefulness of ConveRT-based sentence representations have been recently confirmed on other intent classification datasets (Casanueva et al., 2020), with different intent classifiers (Bunk et al., 2020), and in another dialog task: turn-based value extraction (Coope et al., 2020; Bunk et al., 2020). In future work, we plan to investigate other possible applications of transfer, especially for low-data setups.

5 Conclusion

We have introduced ConveRT, a new light-weight model of neural response selection for dialog, based on Transformer-backed dual-encoder networks, and have demonstrated its state-of-the-art performance on an array of response selection tasks and in transfer learning for intent classification tasks. In addition to offering *more accurate* conversational pretraining models this work has also resulted in *more compact* conversational pretraining. The quantized versions of ConveRT and multi-context ConveRT take up only 59 MB and 73 MB, respectively, and train for 18 hours with a training cost estimate of only 85 USD. In the hope that this work will motivate and guide further developments in the area of retrieval-based task-oriented dialog, we publicly release pretrained ConveRT models.

References

- Rami Al-Rfou, Marc Pickett, Javier Snaider, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2016. Conversational contextual cues: The case of personalization and history for response ranking. *CoRR*, abs/1606.00372.
- Aishwarya Bhandare, Vamsi Sripathi, Deepthi Karkada, Vivek Menon, Sun Choi, Kushal Datta, and Vikram Saletore. 2019. Efficient 8-bit quantization of transformer neural machine language translation model. *CoRR*, abs/1906.00532.
- Basma El Amel Boussaha, Nicolas Hernandez, Christine Jacquin, and Emmanuel Morin. 2019. Deep retrieval-based dialogue systems: A short review. *CoRR*, abs/1907.12878.
- Pawel Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. MultiWOZ - A large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. In *Proceedings of EMNLP*, pages 5016–5026.
- Tanja Bunk, Daksh Varshneya, Vladimir Vlasov, and Alan Nichol. 2020. DIET: Lightweight language understanding for dialogue systems. *CoRR*, abs/2004.09936.
- Iñigo Casanueva, Tadas Temcinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. Efficient intent detection with dual sentence encoders. *CoRR*, abs/2003.04807.
- Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. Universal sentence encoder for English. In *Proceedings of EMNLP*, pages 169–174.
- Debanjan Chaudhuri, Agustinus Kristiadi, Jens Lehmann, and Asja Fischer. 2018. Improving response selection in multi-turn dialogue systems by incorporating domain knowledge. In *Proceedings of CoNLL*, pages 497–507.
- Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang Tang. 2017. A survey on dialogue systems: Recent advances and new frontiers. *CoRR*, abs/1711.01731.
- Muthuraman Chidambaram, Yinfei Yang, Daniel Cer, Steve Yuan, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2019. Learning cross-lingual sentence representations via a multi-task dual-encoder model. In *Proceedings of the 4th Workshop on Representation Learning for NLP*, pages 250–259.
- Sam Coope, Tyler Farghly, Daniela Gerz, Ivan Vulić, and Matthew Henderson. 2020. Span-ConveRT: Few-shot span extraction for dialog with pretrained conversational representations. In *Proceedings of ACL*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186.
- Wenchao Du and Alan Black. 2018. Data augmentation for neural online chats response selection. In *Proceedings of the 2nd International Workshop on Search-Oriented Conversational AI*, pages 52–58.
- Layla El Asri, Hannes Schulz, Shikhar Sharma, Jeremie Zumer, Justin Harris, Emery Fine, Rahul Mehrotra, and Kaheer Suleman. 2017. Frames: A corpus for adding memory to goal-oriented dialogue systems. In *Proceedings of SIGDIAL*, pages 207–219.
- Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, and Dilek Hakkani-Tür. 2019. Multiwoz 2.1: Multi-domain dialogue state corrections and state tracking baselines. *CoRR*, abs/1907.01669.
- Ahmed Fadhil and Gianluca Schiavo. 2019. Designing for health chatbots. *CoRR*, abs/1902.09022.
- Jamie Fraser, Ioannis Papaioannou, and Oliver Lemon. 2018. Spoken conversational AI in video games: Emotional dialogue management increases user engagement. In *Proceedings of IVA*.
- Chulaka Gunasekara, Jonathan K. Kummerfeld, Lazaros Polymenakos, and Walter Lasecki. 2019. DSTC7 task 1: Noetic end-to-end response selection. In *Proceedings of the 1st Workshop on NLP for Conversational AI*, pages 60–67.
- Song Han, Huizi Mao, and William J. Dally. 2016. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. In *Proceedings of ICLR*.
- Charles T. Hemphill, John J. Godfrey, and George R. Doddington. 1990. The ATIS Spoken Language Systems Pilot Corpus. In *Proceedings of the Workshop on Speech and Natural Language*, HLT '90, pages 96–101.
- Matthew Henderson, Rami Al-Rfou, Brian Strope, Yun-Hsuan Sung, László Lukács, Ruiqi Guo, Sanjiv Kumar, Balint Miklos, and Ray Kurzweil. 2017. Efficient natural language response suggestion for smart reply. *CoRR*, abs/1705.00652.
- Matthew Henderson, Pawel Budzianowski, Iñigo Casanueva, Sam Coope, Daniela Gerz, Girish Kumar, Nikola Mrkšić, Georgios Spithourakis, Pei-Hao Su, Ivan Vulić, and Tsung-Hsien Wen. 2019a. A repository of conversational datasets. In *Proceedings of the 1st Workshop on Natural Language Processing for Conversational AI*.
- Matthew Henderson, Blaise Thomson, and Jason D. Williams. 2014. The Second Dialog State Tracking Challenge. In *Proceedings of SIGDIAL*, pages 263–272.

- Matthew Henderson, Ivan Vulić, Daniela Gerz, Iñigo Casanueva, Paweł Budzianowski, Sam Coope, Georgios Spithourakis, Tsung-Hsien Wen, Nikola Mrkšić, and Pei-Hao Su. 2019b. [Training neural response selection for task-oriented dialogue systems](#). In *Proceedings of ACL*, pages 5392–5404.
- Dan Hendrycks and Kevin Gimpel. 2016. [Gaussian error linear units \(GELUs\)](#). *arXiv preprint arXiv:1606.08415*.
- Samuel Humeau, Kurt Shuster, Marie-Anne Lachaux, and Jason Weston. 2020. [Poly-encoders: Transformer architectures and pre-training strategies for fast and accurate multi-sentence scoring](#). In *Proceedings of ICLR*.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. [Albert: A Lite BERT for self-supervised learning of language representations](#). In *Proceedings of ICLR*.
- Liliana Laranjo, Adam G. Dunn, Huong Ly Tong, Ahmet Baki Kocaballi, Jessica Chen, Rabia Bashir, Didi Surian, Blanca Gallego, Farah Magrabi, Annie Y.S. Lau, and Enrico Coiera. 2018. [Conversational agents in healthcare: A systematic review](#). *Journal of the American Medical Informatics Association*, 25(9):1248–1258.
- Bing Liu, Gökhan Tür, Dilek Hakkani-Tür, Pararth Shah, and Larry P. Heck. 2018. [Dialogue learning with human teaching and feedback in end-to-end trainable task-oriented dialogue systems](#). In *Proceedings of NAACL-HLT*, pages 2060–2069.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [RoBERTa: A robustly optimized BERT pretraining approach](#). *CoRR*, abs/1907.11692.
- Ryan Thomas Lowe, Nissan Pow, Iulian Vlad Serban, Laurent Charlin, Chia-Wei Liu, and Joelle Pineau. 2017. [Training end-to-end dialogue systems with the ubuntu dialogue corpus](#). *Dialogue & Discourse*, 8(1):31–65.
- Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. *Introduction to Information Retrieval*. Cambridge University Press.
- Pierre-Emmanuel Mazaré, Samuel Humeau, Martin Raison, and Antoine Bordes. 2018. [Training millions of personalized dialogue agents](#). In *Proceedings of EMNLP*, pages 2775–2779.
- Shikib Mehri, Evgeniia Razumovskaia, Tiancheng Zhao, and Maxine Eskenazi. 2019. [Pretraining methods for dialog context representation learning](#). In *Proceedings of ACL*, pages 3836–3845.
- Paul Michel, Omer Levy, and Graham Neubig. 2019. [Are sixteen heads really better than one?](#) In *Proceedings of NeurIPS*.
- Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory F. Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. 2018. [Mixed precision training](#). In *Proceedings of ICLR*.
- Seyed-Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, and Hassan Ghasemzadeh. 2019. [Improved knowledge distillation via teacher assistant: Bridging the gap between student and teacher](#). *CoRR*, abs/1902.03393.
- Nikola Mrkšić, Diarmuid Ó Séaghdha, Blaise Thomson, Milica Gašić, Pei-Hao Su, David Vandyke, Tsung-Hsien Wen, and Steve Young. 2015. [Multi-domain dialog state tracking using recurrent neural networks](#). In *Proceedings of ACL*, pages 794–799.
- Nikola Mrkšić, Ivan Vulić, Diarmuid Ó Séaghdha, Ira Leviant, Roi Reichart, Milica Gašić, Anna Korhonen, and Steve Young. 2017. [Semantic specialisation of distributional word vector spaces using monolingual and cross-lingual constraints](#). *Transactions of the ACL*, pages 314–325.
- Gabriel Pereyra, George Tucker, Jan Chorowski, Lukasz Kaiser, and Geoffrey E. Hinton. 2017. [Regularizing neural networks by penalizing confident output distributions](#). *CoRR*, abs/1701.06548.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#). *OpenAI Blog*, 1(8).
- Osman Ramadan, Paweł Budzianowski, and Milica Gašić. 2018. [Large-scale multi-domain belief tracking with knowledge sharing](#). In *Proceedings of ACL*, pages 432–437.
- Antoine Raux, Brian Langner, Alan W. Black, and Maxine Eskenazi. 2003. [LET’s GO: Improving spoken dialog systems for the elderly and non-natives](#). In *Proceedings of EUROSPEECH*.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. [Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter](#). *CoRR*, abs/1910.01108.
- Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. 2019. [Green AI](#). *CoRR*, abs/1907.10597.
- Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018. [Self-attention with relative position representations](#). In *Proceedings of NAACL-HLT*, pages 464–468.
- Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W. Mahoney, and Kurt Keutzer. 2019. [Q-BERT: Hessian based ultra low precision quantization of BERT](#). *CoRR*, abs/1909.05840.

- Emma Strubell, Ananya Ganesh, and Andrew McCallum. 2019. [Energy and policy considerations for deep learning in NLP](#). In *Proceedings of ACL*, pages 3645–3650.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2016. [Re-thinking the inception architecture for computer vision](#). In *Proceedings of CVPR*, pages 2818–2826.
- Raphael Tang, Yao Lu, Linqing Liu, Lili Mou, Olga Vechtomova, and Jimmy Lin. 2019. [Distilling task-specific knowledge from BERT into simple neural networks](#). *CoRR*, abs/1903.12136.
- Ashish Vaswani, Samy Bengio, Eugene Brevdo, François Chollet, Aidan Gomez, Stephan Gouws, Llion Jones, Łukasz Kaiser, Nal Kalchbrenner, Niki Parmar, Ryan Sepassi, Noam Shazeer, and Jakob Uszkoreit. 2018. [Tensor2Tensor for neural machine translation](#). In *Proceedings of the 13th Conference of the Association for Machine Translation in the Americas*, pages 193–199. Association for Machine Translation in the Americas.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Proceedings of NeurIPS*, pages 6000–6010.
- Jesse Vig and Kalai Ramea. 2019. [Comparison of transfer-learning approaches for response selection in multi-turn conversations](#). In *Proceedings of DSTC-7*.
- Vladimir Vlasov, Johannes E. M. Mosig, and Alan Nichol. 2019. [Dialogue transformers](#). *CoRR*, abs/1910.00486.
- Mengting Wan and Julian McAuley. 2016. [Modeling ambiguity, subjectivity, and diverging viewpoints in opinion question answering systems](#). In *Proceedings of ICDM*, pages 489–498.
- Hao Wang, Zhengdong Lu, Hang Li, and Enhong Chen. 2013. [A dataset for research on short-text conversations](#). In *Proceedings of EMNLP*, pages 935–945.
- Tsung-Hsien Wen, Yishu Miao, Phil Blunsom, and Steve J. Young. 2017a. [Latent intention dialogue models](#). In *Proceedings of ICML*, pages 3732–3741.
- Tsung-Hsien Wen, David Vandyke, Nikola Mrkšić, Milica Gašić, Lina M. Rojas-Barahona, Pei-Hao Su, Stefan Ultes, and Steve Young. 2017b. [A network-based end-to-end trainable task-oriented dialogue system](#). In *Proceedings of EACL*, pages 438–449.
- Jason Weston, Emily Dinan, and Alexander Miller. 2018. [Retrieve and refine: Improved sequence generation models for dialogue](#). In *Proceedings of the 2018 EMNLP Workshop SCAI: The 2nd International Workshop on Search-Oriented Conversational AI*, pages 87–92.
- Jason Williams. 2012. [A critical analysis of two statistical spoken dialog systems in public use](#). In *Proceedings of SLT*.
- Jason D. Williams. 2014. [Web-style ranking and SLU combination for dialog state tracking](#). In *Proceedings of SIGDIAL*, pages 282–291.
- Yu Wu, Wei Wu, Chen Xing, Ming Zhou, and Zhoujun Li. 2017. [Sequential matching network: A new architecture for multi-turn response selection in retrieval-based chatbots](#). In *Proceedings of ACL*, pages 496–505.
- Yinfei Yang, Steve Yuan, Daniel Cer, Sheng-Yi Kong, Noah Constant, Petr Pilar, Heming Ge, Yun-hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. [Learning semantic textual similarity from conversations](#). In *Proceedings of the 3rd Workshop on Representation Learning for NLP*, pages 164–174.
- Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. 2019. [Q8BERT: Quantized 8bit BERT](#). *CoRR*, abs/1910.06188.
- Matthew D. Zeiler. 2012. [ADADELTA: an adaptive learning rate method](#). *CoRR*, abs/1212.5701.
- Sanqiang Zhao, Raghav Gupta, Yang Song, and Denny Zhou. 2019a. [Extreme language model compression with optimal subwords and shared projections](#). *CoRR*, abs/1909.11687.
- Tiancheng Zhao, Kaige Xie, and Maxine Eskénazi. 2019b. [Rethinking action spaces for reinforcement learning in end-to-end dialog agents with latent variable models](#). In *Proceedings of NAACL-HLT*, pages 1208–1218.
- Xiangyang Zhou, Lu Li, Daxiang Dong, Yi Liu, Ying Chen, Wayne Xin Zhao, Dianhai Yu, and Hua Wu. 2018. [Multi-turn response selection for chatbots with deep attention matching network](#). In *Proceedings of ACL*, pages 1118–1127.