

Progressive Memory Banks for Incremental Domain Adaptation

Nabiha Asghar

Collaborators:

Lili Mou, Kira Selby, Kevin Pantasdo, Pascal Poupart, Xin Jiang



Motivation

- **Domain Adaptation (DA):** Transfer knowledge from one domain to another (in a machine learning system; especially neural networks)
- **Incremental Domain Adaptation (IDA):** Sequentially incoming domains
 - Only have access to data of current domain
 - Build a unified model that performs well on all domains
- **Use-cases of IDA**
 1. Company loses a client and its data, but wants to preserve the 'knowledge' in the ML system
 2. Quickly adapt to new domain/data without training from scratch
 3. Don't know the domain of a data point during inference

Outline

- Prevalent and State-of-the-art DA & IDA methods in NLP
- Proposed Approach: Progressive Memory for IDA
- Theoretical Analysis
- Empirical Experiments
 - Natural Language Inference (Classification)
 - Dialogue Response Generation
- Conclusion

Related Work - DA & IDA

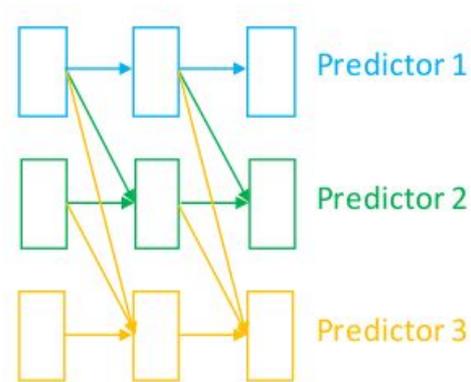
- **Multi-task learning:** Jointly train on all domains
 - Non-incremental DA
 - Expensive to add new domain; needs data for all domains

Related Work - DA & IDA

- **Multi-task learning:** Jointly train on all domains
 - Non-incremental DA
 - Expensive to add new domain; needs data for all domains
- **Finetuning:** Sequentially train on all domains
 - Catastrophic forgetting of old domains

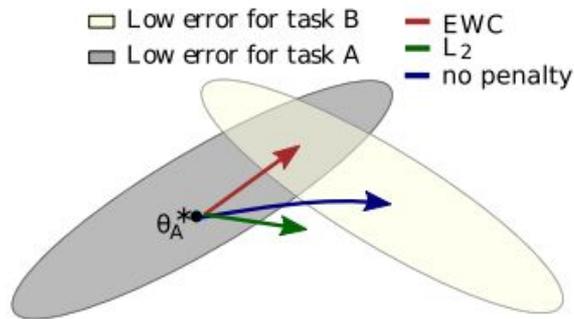
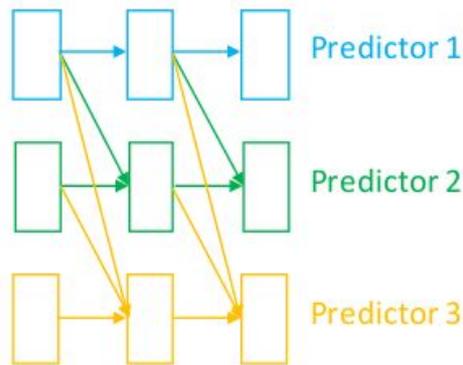
Related Work - DA & IDA

- **Multi-task learning:** Jointly train on all domains
 - Non-incremental DA
 - Expensive to add new domain; needs data for all domains
- **Finetuning:** Sequentially train on all domains
 - Catastrophic forgetting of old domains
- **Progressive Neural Networks:** Training with network expansion and partial freezing
 - For prediction, need to know domain of input



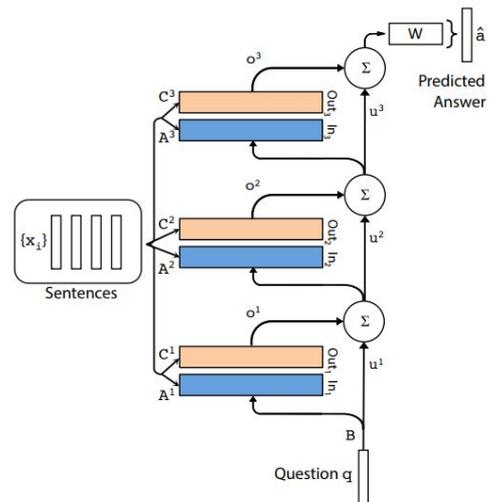
Related Work - DA & IDA

- **Multi-task learning:** Jointly train on all domains
 - Non-incremental DA
 - Expensive to add new domain; needs data for all domains
- **Finetuning:** Sequentially train on all domains
 - Catastrophic forgetting of old domains
- **Progressive Neural Networks:** Training with network expansion and partial freezing
 - For prediction, need to know domain of input
- **Elastic Weight Consolidation (EWC):** Finetuning with regularization
 - Control learning on weights important for older domains
 - keeps the weights in a neighborhood of one possible minimizer of the empirical risk of the first task
 - needs to store a large number of parameters

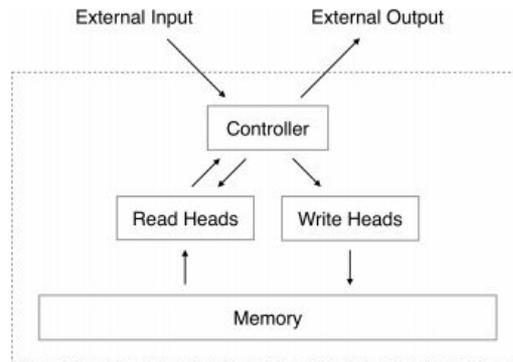


Related Work - Memory Networks

- **End-to-end memory network**
 - Assign a memory slot to an input sentence/sample
 - Assign a memory slot to one history

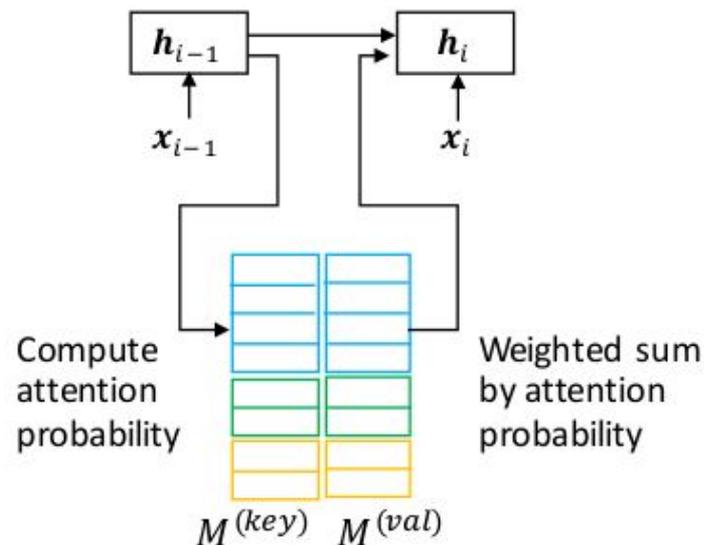


- **Neural Turing Machine**
 - Memory is not directly parameterized; read/written by neural controller
 - Serves as temporary scratch paper; does not store knowledge



Proposed Approach - Progressive Memory

- Incrementally increase model capacity (by increasing memory size)
- Memory slots store knowledge in distributed fashion
- We adopt key-value memory



Progressive Memory

At time step i :

The RNN state is given by $h_i = \text{RNN}(h_{i-1}, x_i)$

The memory mechanism computes an attention probability α_i by

$$\tilde{\alpha}_{i,j} = \exp\{h_{i-1}^\top m_j^{(\text{key})}\}$$

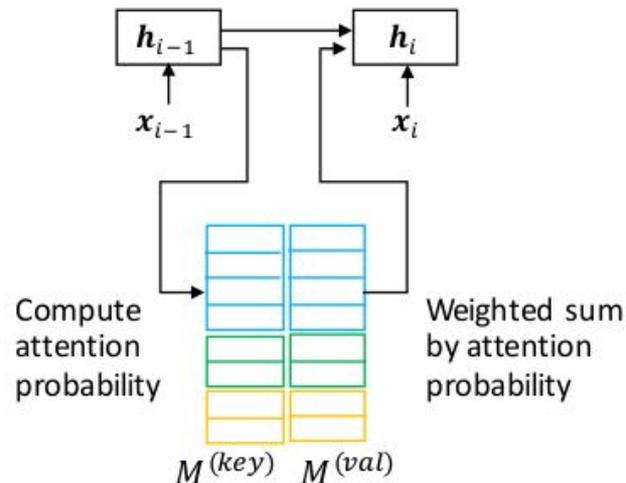
$$\alpha_{i,j} = \frac{\tilde{\alpha}_{i,j}}{\sum_{j'=1}^N \tilde{\alpha}_{i,j'}}$$

$m_j^{(\text{key})}$: key vector of j 'th memory slot (N in total)

Retrieve memory content by weighted sum (by attention probability) of all memory values:

$$c_i = \sum_{j=1}^N \alpha_{i,j} m_j^{(\text{val})}$$

$m_j^{(\text{val})}$: value vector of j 'th memory slot



$$h_i = \text{RNN}(h_{i-1}, [x_i, c_i])$$

Progressive Memory

At time step i :

The RNN state is given by $\mathbf{h}_i = \text{RNN}(\mathbf{h}_{i-1}, \mathbf{x}_i)$

The memory mechanism computes an attention probability α_i by

$$\tilde{\alpha}_{i,j} = \exp\{\mathbf{h}_{i-1}^\top \mathbf{m}_j^{(\text{key})}\}$$
$$\alpha_{i,j} = \frac{\tilde{\alpha}_{i,j}}{\sum_{j'=1}^N \tilde{\alpha}_{i,j'}}$$

$\mathbf{m}_j^{(\text{key})}$: key vector of j 'th memory slot (N in total)

Retrieve memory content by weighted sum (by attention probability) of all memory values:

$$\mathbf{c}_i = \sum_{j=1}^N \alpha_{i,j} \mathbf{m}_j^{(\text{val})}$$

$\mathbf{m}_j^{(\text{val})}$: value vector of j 'th memory slot

For IDA:

Add M slots to original N slots

$$\alpha_{i,j}^{(\text{expand})} = \frac{\tilde{\alpha}_{i,j}}{\sum_{j'=1}^{N+M} \tilde{\alpha}_{i,j'}}$$
$$\mathbf{c}_i^{(\text{expand})} = \sum_{j=1}^{N+M} \alpha_{i,j}^{(\text{expand})} \mathbf{m}_j^{(\text{val})}$$

Algorithm

Algorithm 1: Progressive Memory for IDA

Input: A sequence of domains D_0, D_1, \dots, D_n

Output: A model performing well on all domains

Initialize a memory-augmented RNN

Train the model on D_0

for D_1, \dots, D_n **do**

 Expand the memory with new slots

 Load RNN weights and existing memory banks

 Train the model by updating all parameters

end

Return: The resulting model

Training Considerations

- Freezing learned params **versus** Finetuning learned params
 - Empirical results are better for latter
- Finetuning w/o increasing memory **versus** Finetuning w/ increasing memory
 - increased model capacity helps to learn new domain with less overriding of the previously learned model. Empirical results confirm this.
- Expanding hidden states **versus** Expanding Memory
 - An alternate way of increasing model capacity
 - similar to the progressive neural network, except that all weights are fine-tuned and there are connections from new states to existing states.
 - Theoretical and empirical results show latter is better

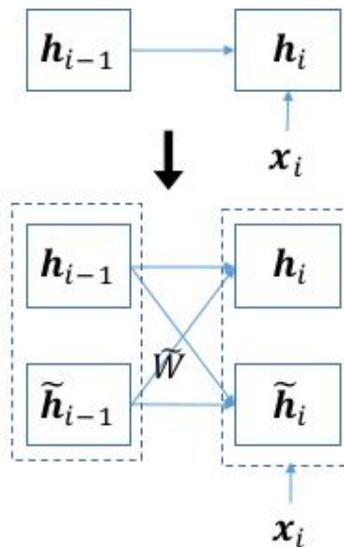
Expanding hidden states vs Expanding Memory

Theorem 1. *Let RNN have vanilla transition with the linear activation function, and let the RNN state at the last step \mathbf{h}_{i-1} be fixed. For a particular data point, if the memory attention satisfies $\sum_{j=N+1}^{N+M} \tilde{\alpha}_{i,j} \leq \sum_{j=1}^N \tilde{\alpha}_{i,j}$, then memory expansion yields a lower expected mean squared difference in \mathbf{h}_i than RNN state expansion, under reasonable assumptions. That is,*

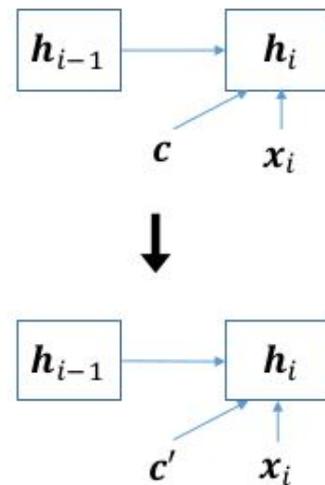
$$\mathbb{E} \left[\|\mathbf{h}_i^{(m)} - \mathbf{h}_i\|^2 \right] \leq \mathbb{E} \left[\|\mathbf{h}_i^{(s)} - \mathbf{h}_i\|^2 \right]$$

where $\mathbf{h}_i^{(m)}$ refers to the hidden states if the memory is expanded. $\mathbf{h}_i^{(s)}$ refers to the original dimensions of the RNN states, if we expand the size of RNN states themselves.

(a) Expand RNN states



(b) Expand memory

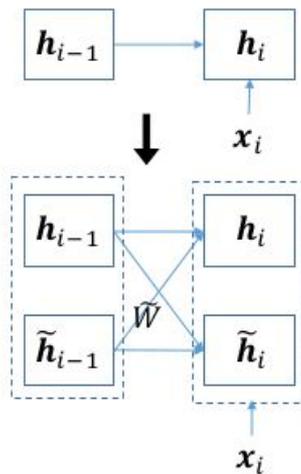


To prove: $\mathbb{E} [\|\mathbf{h}_i^{(m)} - \mathbf{h}_i\|^2] \leq \mathbb{E} [\|\mathbf{h}_i^{(s)} - \mathbf{h}_i\|^2]$

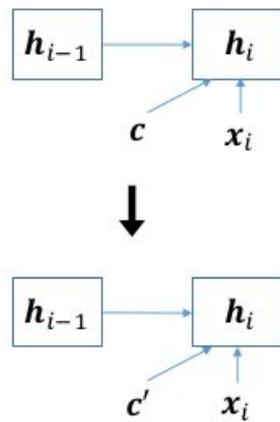
Suppose the original hidden state \mathbf{h}_i is D -dimensional. We assume each memory slot is d -dimensional, and that the additional RNN units when expanding the hidden state are also d -dimensional. We further assume every variable in the expanded memory and expanded weights are iid with zero mean and variance σ^2 . Finally, every variable in the learned memory slots, i.e., m_{jk} , follows the same distribution (zero mean, variance σ^2). This assumption may not be true after the network is trained, but is useful for proving theorems.

$$\begin{aligned}
 & \mathbb{E} [\|\mathbf{h}_i^{(s)} - \mathbf{h}_i\|^2] \\
 &= \mathbb{E} [\|\widetilde{W} \cdot \widetilde{\mathbf{h}}_{i-1}\|^2] \\
 &= \sum_{j=1}^D \sum_{i=1}^d \mathbb{E} [(\widetilde{w}_{jk})^2] \mathbb{E} [(\widetilde{h}_{i-1}[k])^2] \\
 &= D \cdot d \cdot \text{Var}(w) \cdot \text{Var}(h) \\
 &= Dd\sigma^2\sigma^2
 \end{aligned}$$

(a) Expand RNN states



(b) Expand memory



$$\begin{aligned}
 & \mathbb{E} [\|\mathbf{h}_i^{(m)} - \mathbf{h}_i\|^2] \\
 &= \mathbb{E} [\|W_{(c)} \Delta \mathbf{c}\|^2] \text{ where } \Delta \mathbf{c} \stackrel{\text{def}}{=} \mathbf{c}' - \mathbf{c} \\
 &= Dd\sigma^2 \text{Var}(\Delta c_k)
 \end{aligned}$$

$W_{(c)}$ is the weight matrix connecting attention content to RNN states.

It remains to show that $\text{Var}(\Delta c_k) \leq \sigma^2$

$$\Delta \mathbf{c} = \mathbf{c}' - \mathbf{c}$$

$$= \sum_{j=1}^N (\alpha'_j - \alpha_j) \mathbf{m}_j + \sum_{j=N+1}^{N+M} \alpha'_j \mathbf{m}_j = \sum_{j=1}^{N+M} \beta_j \mathbf{m}_j$$

$$\beta_j \stackrel{\text{def}}{=} \begin{cases} \frac{-\tilde{\alpha}_j \frac{\tilde{\alpha}_{N+1} + \dots + \tilde{\alpha}_{N+M}}{\tilde{\alpha}_1 + \dots + \tilde{\alpha}_N}}{\tilde{\alpha}_1 + \dots + \tilde{\alpha}_{N+M}}, & \text{if } 1 \leq j \leq N \\ \frac{\tilde{\alpha}_j}{\tilde{\alpha}_1 + \dots + \tilde{\alpha}_{N+M}}, & \text{if } N+1 \leq j \leq N+M \end{cases}$$

$$\text{Var}(\Delta c_k) = \mathbb{E}[(c'_k - c_k)^2] \quad \forall 1 \leq k \leq d$$

$$= \frac{1}{d} \mathbb{E} [\|\mathbf{c}' - \mathbf{c}\|^2]$$

$$= \frac{1}{d} \mathbb{E} \left[\sum_{k=1}^d \left(\sum_{j=1}^{N+M} \beta_j m_{jk} \right)^2 \right]$$

$$\leq \sigma^2 \mathbb{E} \left[\sum_{j=1}^{N+M} (\alpha'_j)^2 \right]$$

$$\leq \sigma^2$$

□

Memory	Unnormalized measure	Original attn. prob.	Expanded attn. prob.
\mathbf{m}_1	$\tilde{\alpha}_1$	α_1	α'_1
\mathbf{m}_2	$\tilde{\alpha}_2$	α_2	α'_2
...
\mathbf{m}_N	$\tilde{\alpha}_N$	α_N	α'_N
\mathbf{m}_{N+1}	$\tilde{\alpha}_{N+1}$		α'_{N+1}
...
\mathbf{m}_{N+M}	$\tilde{\alpha}_{N+M}$		α'_{N+M}

Figure 3: Attention probabilities before and after memory expansion.

Competing Methods

- Multi-task learning (Non-IDA)
- Finetuning with fixed memory
- **Finetuning with increasing memory**
- Finetuning with expanding hidden states
- Progressive Neural Network
- Elastic Weight Consolidation (EWC)

Competing Methods

- Multi-task learning (Non-IDA)
- Finetuning with fixed memory*
- **Finetuning with increasing memory***
- Finetuning with expanding hidden states*
- Progressive Neural Network
- Elastic Weight Consolidation (EWC)

* with and without additional vocabulary

Experiment I - Natural Language Inference

- Classification Task
 - Determine the relationship between two sentences (*entailment*, *contradiction* or *neutral*)
- Dataset: MultiNLI Corpus (~400K labelled samples)
 - 5 domains: Fiction, Government, Slate, Telephone, Travel
- Base Model: BiLSTM network with pretrained GloVe embeddings

Experiment I - Results

#Line	Model	Trained on/by	% Accuracy on	
			S	T
1	RNN	S	65.01 \downarrow	61.23 \downarrow
2		T	56.46 \downarrow	66.49 \downarrow
3	RNN+ Mem	S	65.41 \downarrow	60.87 \downarrow
4		T	56.77 \downarrow	67.01 \downarrow
5		S+T	66.02 \downarrow	70.00
6	RNN + Mem	S \rightarrow T (F)	65.62 \downarrow	69.90 \downarrow
7		S \rightarrow T (F+M)	66.23	70.21
8		S \rightarrow T (F+M+V)	67.55	70.82
9		S \rightarrow T (F+H)	64.09 \downarrow	68.35 \downarrow
10		S \rightarrow T (F+H+V)	63.68 \downarrow	68.02 \downarrow
11		S \rightarrow T (EWC)	66.02 \downarrow	64.10 \downarrow
12		S \rightarrow T (Progressive)	64.47 \downarrow	68.25 \downarrow

For the statistical test (compared with Line 8), \uparrow, \downarrow : $p < 0.05$ and \uparrow, \downarrow : $p < 0.01$.

Experiment I - Results

Group	Setting	Fic	Gov	Slate	Tel	Travel
Non-IDA	In-domain training	65.41 \downarrow	67.01 \downarrow	59.30 \downarrow	67.20 \downarrow	64.70 \downarrow
	Fic + Gov + Slate + Tel + Travel (multi-task)	70.60\uparrow	73.30	63.80	69.15	67.07 \downarrow
IDA	Fic \rightarrow Gov \rightarrow Slate \rightarrow Tel \rightarrow Travel (F+V)	67.24 \downarrow	70.82 \downarrow	62.41 \downarrow	67.62 \downarrow	68.39
	Fic \rightarrow Gov \rightarrow Slate \rightarrow Tel \rightarrow Travel (F+V+M)	69.36	72.47	63.96	69.74	68.39
	Fic \rightarrow Gov \rightarrow Slate \rightarrow Tel \rightarrow Travel (EWC)	67.12 \downarrow	68.71 \downarrow	59.90 \downarrow	66.09 \downarrow	65.70 \downarrow
	Fic \rightarrow Gov \rightarrow Slate \rightarrow Tel \rightarrow Travel (Progressive)	65.22 \downarrow	67.87 \downarrow	61.13 \downarrow	66.96 \downarrow	67.90

Experiment II - Dialogue Response Generation

- Generation Task
 - Given an input sentence, generate an appropriate output sentence
- Datasets
 - Source Domain: Cornell Movie Corpus (~220K labelled samples)
 - Target Domain: Ubuntu Dialogue Corpus (~15K labelled samples)
- Base Model: Seq2Seq with decoder-to-encoder-attention

Experiment II - Results

# Line	Model	Trained on/by	BLEU-2 on		W2V-Sim on	
			S	T	S	T
1	RNN	S	2.842 \uparrow	0.738 \downarrow	0.480 \downarrow	0.456 \downarrow
2		T	0.795 \downarrow	1.265 \downarrow	0.454 \downarrow	0.480 \downarrow
3	RNN+ Mem	S	3.074 \uparrow	0.712 \downarrow	0.498 \downarrow	0.471 \downarrow
4		T	0.920 \downarrow	1.287 \downarrow	0.462 \downarrow	0.487 \downarrow
5		S+T	2.650 \uparrow	0.889 \downarrow	0.471 \downarrow	0.462 \downarrow
6	RNN + Mem	S \rightarrow T (F)	1.210 \downarrow	1.101 \downarrow	0.509 \downarrow	0.514 \downarrow
7		S \rightarrow T (F+M)	1.435 \downarrow	1.207 \downarrow	0.526	0.522
8		S \rightarrow T (F+M+V)	1.637	1.652	0.522	0.525
9		S \rightarrow T (F+H)	1.036 \downarrow	1.606 \downarrow	0.503 \downarrow	0.495 \downarrow
10		S \rightarrow T (F+H+V)	1.257 \downarrow	1.419 \downarrow	0.504 \downarrow	0.492 \downarrow
11		S \rightarrow T (EWC)	1.397 \downarrow	1.382 \downarrow	0.513 \downarrow	0.514 \downarrow
12		S \rightarrow T (Progressive)	1.299 \downarrow	1.408 \downarrow	0.502 \downarrow	0.503 \downarrow

Conclusion

- Proposed progressive memory for IDA (Incremental Domain Adaptation)
 - Outperforms other IDA approaches
 - Empirical results show it avoids catastrophic forgetting
 - Theoretical results show it is better than other ways of capacity expansion
-
- Details: <https://arxiv.org/pdf/1811.00239.pdf>

References

- **EWC:** Kirkpatrick et al. “Overcoming catastrophic forgetting in neural networks”. PNAS, 2017.
- **Progressive Neural Networks:** Rusu et al. “Progressive neural networks”. arXiv:1606.04671, 2016.
- **End-to-end Memory Networks:** Sukhbaatar et al. “End-to-end memory networks”. NIPS, 2015.
- **Neural Turing Machine:** Graves et al. “Hybrid computing using a neural network with dynamic external memory”. Nature, 2016.
- **MultiNLI corpus:** Williams et al. “A broad-coverage challenge corpus for sentence understanding through inference”. NAACL, 2018.
- **Cornell Movie Corpus:** Danescu and Lee. “Chameleons in imagined conversations: A new approach to understanding coordination of linguistic style in dialogs”. CMCL, 2011.
- **Ubuntu Dialogue Corpus:** Lowe et al. “The Ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems”. SIGDIAL, 2015.