

Knowledge-enhanced NLP: Progress and Challenge

Nan Duan (段楠)

Lead Researcher

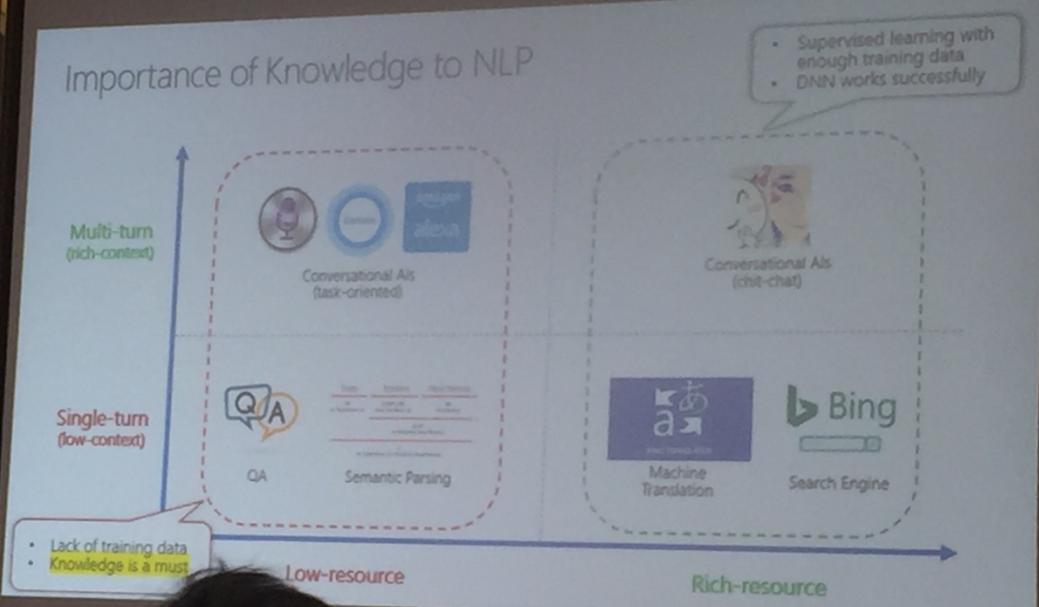
Natural Language Computing Group

Microsoft Research Asia

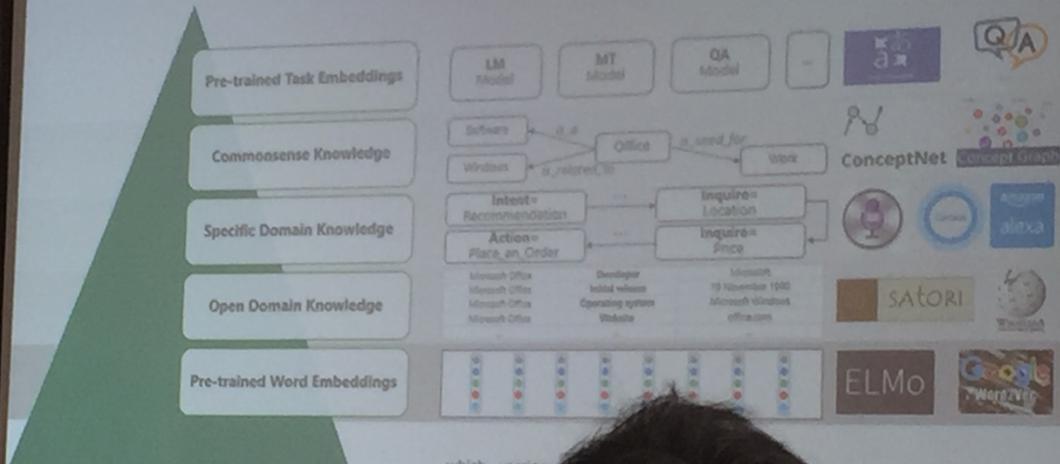
2018-08-30 @ NLPCC 2018 Workshop on QA in NLP



Importance of Knowledge to NLP

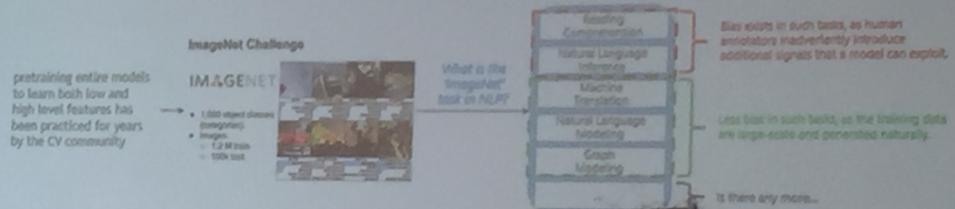


Knowledge Pyramid in NLP

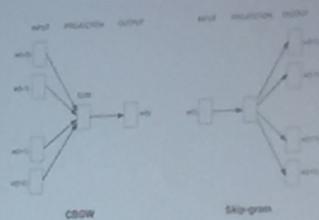


Pre-training: From ImageNet (CV) to ELMo (NLP)

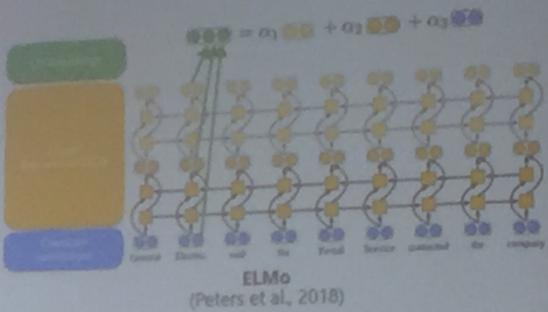
- Shallow pre-training in NLP (such as **Word2Vec** and **GloVe**)
 - Initialize the first layer of a neural network
 - This is like initializing a CV model with pretrained representations that only encode edges
- Deep pre-training in NLP (such as **ELMo** and **OpenAI Transformer**)
 - Pretrain the entire model with hierarchical representations
 - This is like learning the full hierarchy of CV features, from edges to shapes to high-level semantic concepts



Pre-trained Word Embeddings



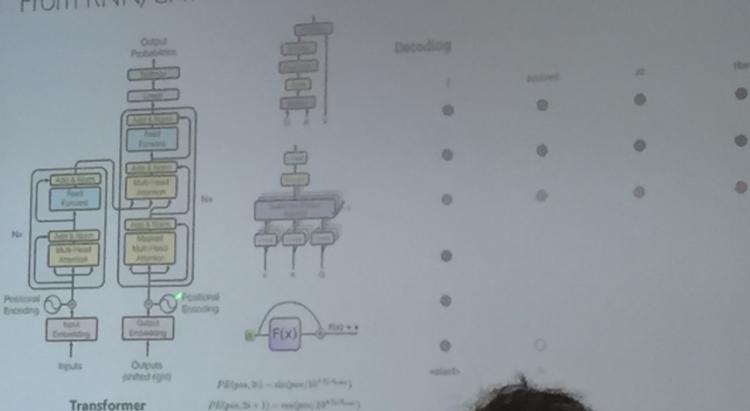
Word2Vec
(Mikolov et al., 2013)



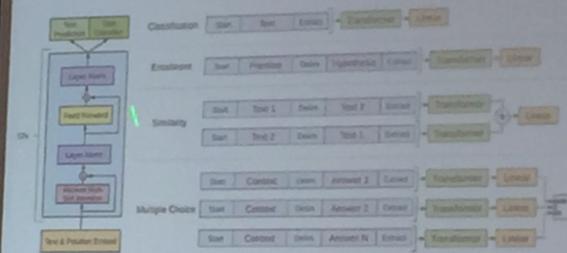
ELMo
(Peters et al., 2018)

TASK	PREVIOUS SOTA	OUR	ELMo + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Lin et al. (2017)	84.4	85.3	4.7 / 2.9%
SNLI	Chen et al. (2017)	88.6	88.7 \pm 0.17	0.7 / 5.8%
SRL	Hermann et al. (2017)	81.7	81.4	3.2 / 7.2%
CoNLL	Dozat et al. (2017)	67.2	67.2	3.2 / 9.8%
NER	Lee et al. (2017)	90.15	90.22 \pm 0.10	2.06 / 21%
STPC	Wang et al. (2017)	51.4	54.7 \pm 0.5	3.3 / 6.8%

From RNN/CNN to Transformer (self-attention)



Unsupervised Pre-training + Supervised Fine-tuning



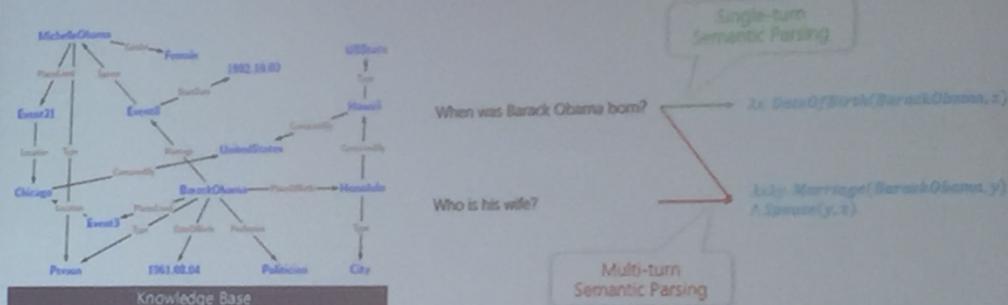
	MNLI-m	MNLI-mm	SST-2	SST-5	QQP	RTE
BERT4LM + ELMo + Attn [6]	-	69.7	79.0	88.2	-	-
ELMo + ELMo + Attn [6]	69.0	80.1	-	-	-	-
Subword Answer Network [23] (ours)	70.0	-	-	-	-	-
ELMo [24]	74.1	75.8	88.3	83.3	-	-
ELMo [24]	75.4	77.5	-	-	82.1	80.2
ELMo + ELMo + Attn [6]	72.2	73.3	-	-	82.7	-
Fine-tuned Transformer 1.3M (ours)	82.1	82.4	89.9	89.3	90.0	-

Method	Story Cloze	RACE-m	RACE-h	RACE
Self-GLS-clip [33]	56.5	-	-	-
Robust Cross-Task Model [7]	77.9	-	-	-
Dynamic Finetune Net [27] (ours)	-	85.4	89.4	81.3
ELMo+BERT [19] (ours)	-	80.7	89.3	83.3
Fine-tuned Transformer 1.3M (ours)	86.8	82.9	87.4	88.0

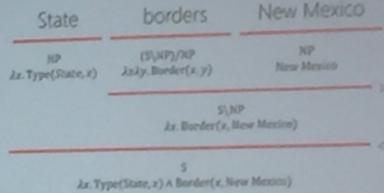
Method	Classification	Semantic Similarity	GLUE			
	CnLP SST2 MRPC	STS2 STS3 QQP RTE	CnLP SST2 MRPC			
Baseline (no pre-train) [6]	-	83.2	-			
ELMo [24] (ours)	-	88.0	-			
TF-GELD [25]	-	-	-			
ELMo+BERT (ours) [19]	-	-	-			
Single-task BERTLM + ELMo + Attn [6]	23.0	90.2	55.3	66.8		
Multi-task BERTLM + ELMo + Attn [6]	33.9	91.6	83.3	72.8	63.3	68.9
Fine-tuned Transformer 1.3M (ours)	48.4	91.3	82.3	80.0	76.3	72.8

Take Semantic Parsing as an Example

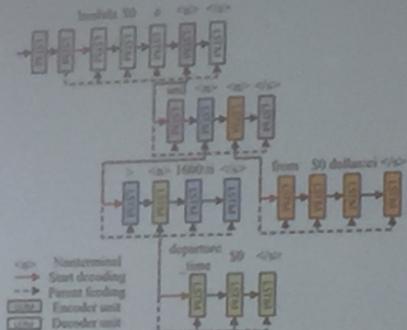
- Convert NL utterances into machine executable LFs based on knowledge base



From Grammars to Neural Networks

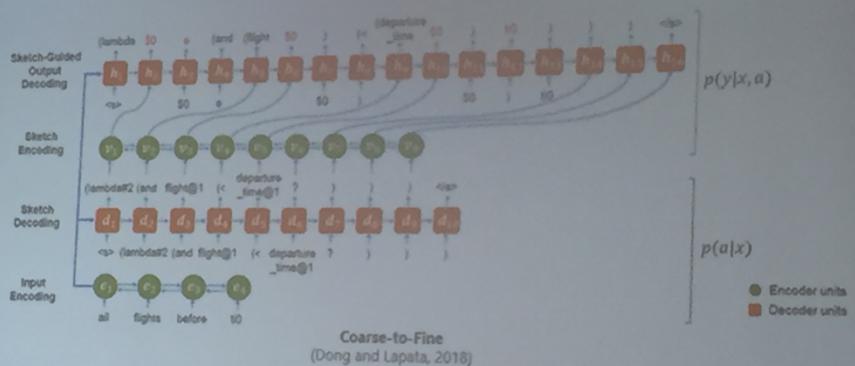


Grammar-based Approach
(Zettlemoyer and Collins, 2007)



Neural Network-based Approach
(Dong and Lapata, 2016)

From One-pass Decoding to Multi-pass Decoding



From Simple Questions to Complex Questions

The SimpleQuestions dataset

This section presents SimpleQuestions, a dataset collected by researchers at University College London using multi-human generated questions. Details and further results on this dataset can be found in the paper: Antoine Bordes, Nicolas Usunier, Sumit Chopra and Jason Weston. Large-Scale Question Answering with Memory Networks, arXiv:1506.02073.

The SimpleQuestions dataset consists of a total of 100,000 questions written in natural language by native English-speaking encyclopedists each paired with a corresponding fact formulated as a triple, mentioning an object, that provides the answer but gives a complete explanation. Factual knowledge is derived from the Knowledge Base Wikipedia. We randomly shuffle these questions and use 20% of them (20,000) as training set, 10% as validation set (10,000), and the remaining 70% as test set.

Here are some examples of questions and facts:

• What previous experience did the author of the book have?

• Who's wife?

• Who's father?

• Who's mother?

• Who's son?

• Who's daughter?

• Who's sibling?

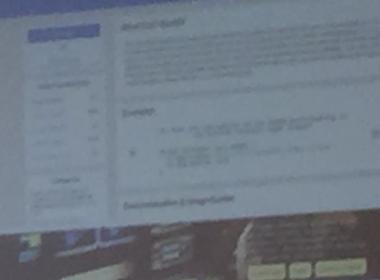
SimpleQuestions Nearly Solved: A New Upperbound and Baseline Approach

Michael Petrukhin
University of Washington Department
of Computer Science & Engineering
mep@cs.washington.edu

Lukas Zettlemoyer
University of Washington Department
of Computer Science & Engineering
lukasz@cs.washington.edu

Because of data ambiguity, that upper bound performance on
this benchmark at **83.4%** (Petrukhin and Zettlemoyer, 2018).

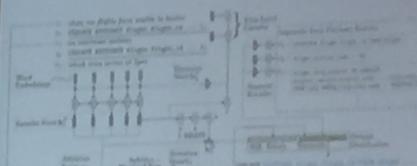
LC-QNAU



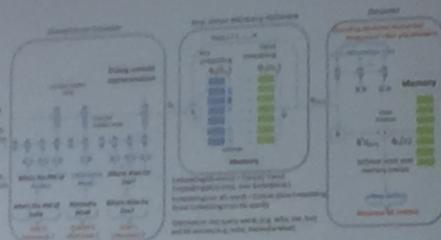
3. Adapted for answering complex
questions that involve reasoning over
multiple facts or concepts

4. Adapted for answering complex
questions that involve reasoning over
multiple facts or concepts

From Single-turn to Multi-turn

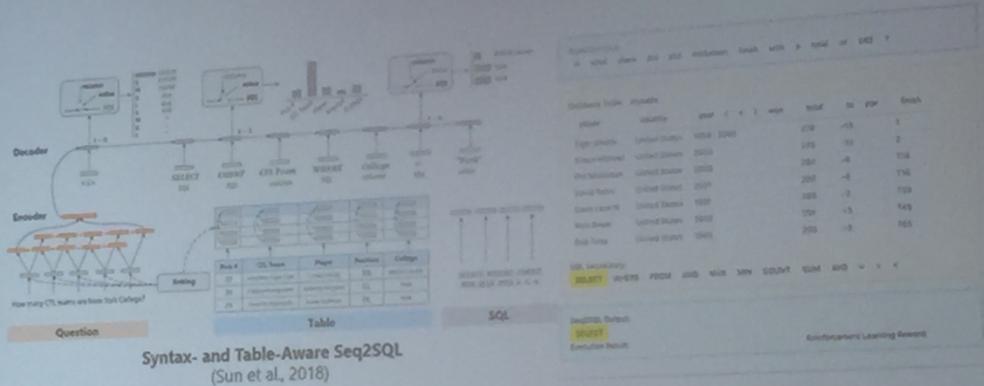


Learning to Map Context-Dependent Sentences to LFs
(Suhr et al., 2018)



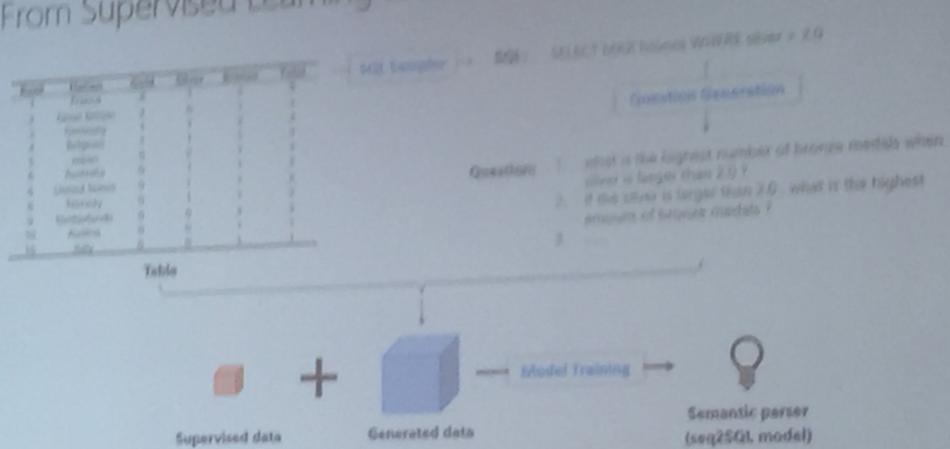
Complex Sequential Question Answering
(Saha et al., 2018)

From Knowledge Graph to Web Tables



Syntax- and Table-Aware Seq2SQL
(Sun et al., 2018)

From Supervised Learning to Semi-supervised Learning

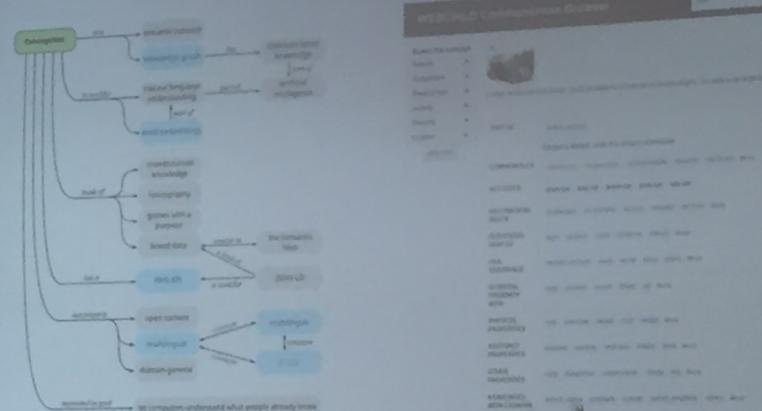


QG Improves Seq2SQL
(Guo et al., 2016)

Conclusion (2)

- Training data is critical to natural language understanding
- Previous work focus on single-turn, single-relation and open domain tasks
- Current work moves to multi-turn, complex-relation and specific tasks
- Take-aways
 - Single-turn questions are almost solved, if we have enough data
 - Multi-turn & complex questions are still challenging

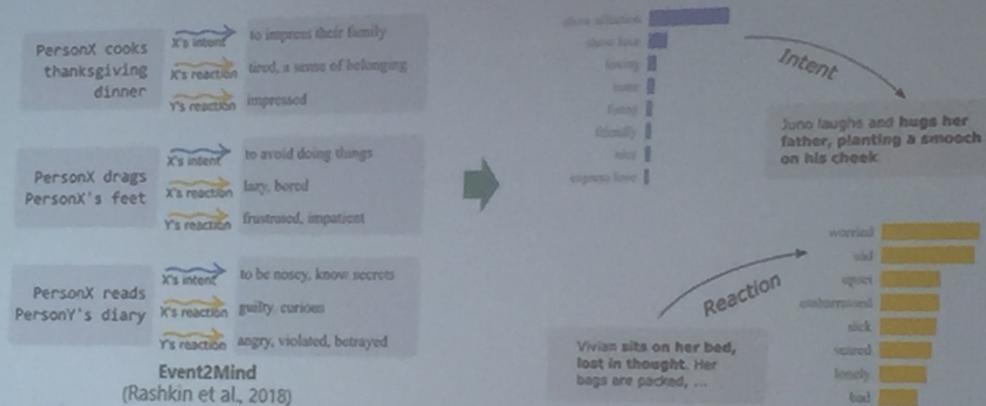
Commonsense Knowledge



ConceptNet

Intent and Reaction Prediction with Commonsense Knowledge

- Use commonsense knowledge learnt from **supervised data**



QA and Reasoning with Commonsense Knowledge

The service was poor, but the food was ...

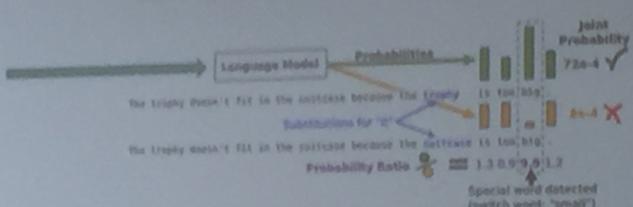
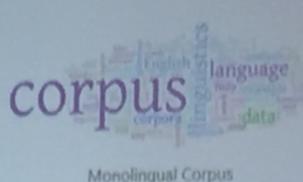
- Use commonsense knowledge learnt from **unsupervised data**

Winograd Schema Challenge

The trophy doesn't fit in the suitcase because it is too big.

Question: What is too big?

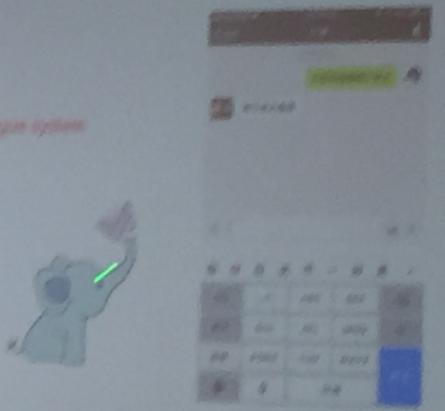
Answer: (a) the trophy (b) the suitcase



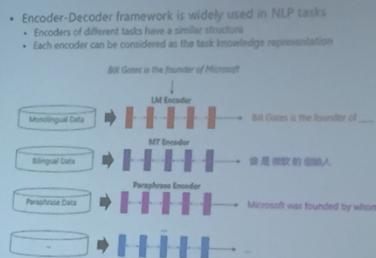
Language Model for Commonsense Reasoning
(Trinh and Le, 2018)

Conclusion (3)

- Commonsense is important to AI models, esp. chatbot and dialogue system
- Few datasets for model developments
- Take-ways
 - Commonsense is invisible most of time
 - But it is important to conversational AIs, such as chatbot and dialogue system



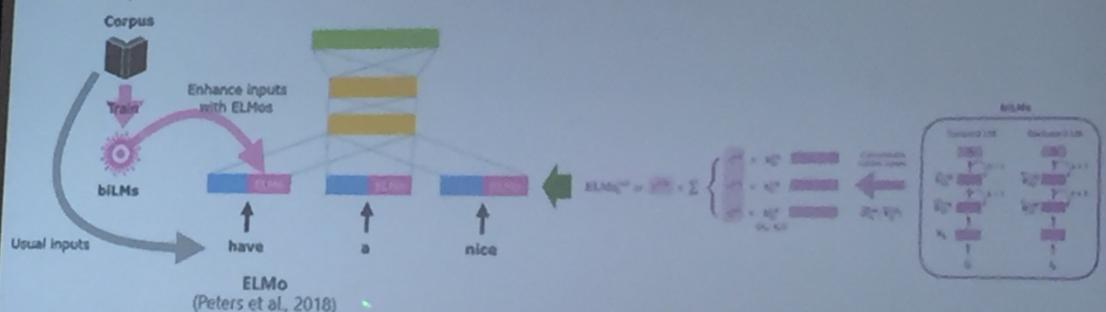
Pre-trained Task Embeddings



Transfer Learning with Pre-trained Task Embeddings in NLP

- Knowledge acquired by a trained model can be transferred during training process for a new task

ELMo can be integrated to almost all neural NLP tasks with simple concatenation to the embedding layer



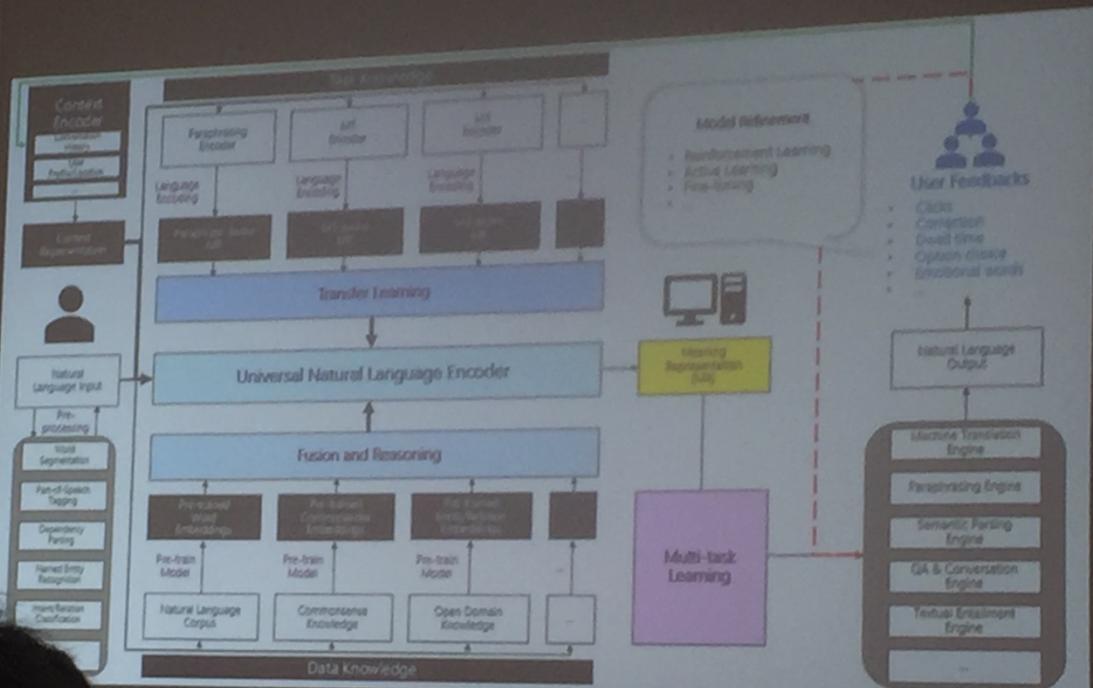
Transfer Learning with Pre-trained Task Embeddings in CV



(Zamir et al., 2018): The task taxonomy (Taskonomy) is a computationally found directed hypergraph that captures the notion of task transferability over any given task dictionary.

Conclusion (4)

- Take-aways
 - Transfer learning is important to NLP



Summary and Challenge

- Models can learn **implicit knowledge** from large-scale datasets for rich-resource tasks!
- Models need **explicit knowledge** for low-resource tasks without much training data!
- How to **define, acquire and represent** knowledge?
- How to **represent meanings of language** with knowledge?
- How to **reason** with knowledge?
- How to **transfer knowledge** from one task/modal to another task/modal?
- How to **memorize, update and forget** knowledge as context?

段楠, 周明. 《智能问答》. 高等教育出版社, 2018 (in press)

