

Lecture 14: Transformers and Self-Attention

Ashish Vaswani and Anna Huang

Learning Representation of Variable Length Data

RNNs

- LSTMs, GRUs and variants dominate recurrent models
- Natural fit for sentences and sequences of pixel

But

Sequential computation inhibits parallelization

No explicit modeling of long and short range dependencies

Attention \Rightarrow why not use attention for representation

Text generation

Classification & regression with self attention

[Parikh et al. 2016, Lin et al. 2016]

Self attention with RNNs [Long et al. 2016, Shao, Grews et al. 2017]

Recurrent attention [Sukhbaatar et al. 2015]

Attention is cheap

[FLOPS]

Self Attention $O(\text{length}^2 \cdot \text{dim})$ if $\text{length} \leq \text{dim}$

RNN(LSTM) $O(\text{length} \cdot \text{dim}^2)$

Convolution $O(\text{length} \cdot \text{dim}^2 \cdot \text{kernel-width})$

@ Multi-head attention:

this is to put attention in different positions
separately

Importance of Residuals

Residuals carry positional information to higher layer, among other information

Music Generation Finding similar motifs

Probabilistic Image Generation

Texture Synthesis with Self Similarity

Non-local NN (Wang 2015)

Selfattention

$$A(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right)V$$

Image Transformer

Combining Locality with Self Attention

Image Transformer (ICML 2018)

Music generation using Self Attention

Music Transformer (ICLR 2019)

→ score → performance sound
composer performance instrument → listener

$$\text{MultiAttention} = \text{softmax} \left(QK^T + \text{score}(QE_{\text{rel}}^T) \right)$$

Equivalence

Relative positions:

Translational Equivalence

- Fast Decoding sequence models ICML 2018 Kaiser

- Better understanding of Vector Quantized auto encoders Roy 2018

- Blockwise Parallel Decoding for Autoregressive model (NeurIPS, 2019) Stern

Lecture 15: Natural Language Generation (NLG)

Neural approaches to NLG

Section: LM and decoding algs: NLG

↳ sub component

- Machine Translation
- (Abstractive) summarization
- Dialogue (chit-chat and task-based)

- Creative writing: storytelling, poetry generation

- free-form QA (i.e. answer is generated)

- Image Captioning

Language Modeling: The task of predicting the next word, given the words so far

$$P(y_t | y_1, y_2, \dots, y_{t-1})$$

- A system that produces this probability distribution is called a Language Model

Conditional Language model:

$$P(y_t | y_1, y_2, \dots, y_{t-1} | x)$$

Decoding algo

- Greedy Decoding

- Beam Search

→ search algorithm which aims to find a high probability sequence (non-necessary the optimal sequence, though) by tracking multiple possible sequences at once

what's the effect of changing beam size k?

• small k has similar problems to greedy decoding

→ ungrammatical, unnatural, nonsensical, incorrect

• Larger k means you consider more hypotheses

→ increasing k reduces some of the problem

above

→ Larger k is more computationally expensive

→ Increasing k too much decreases BLEU

score. This is primarily because large-k

beam search produces too short translation

→ In open ended tasks, large-k can make it output
more generic

Sampling-based decoding

• Pure sampling

→ On each step t, randomly sample from

the prob dist P_t to obtain your next word

→ Like greedy decoding but sample based
instead of argmax

Top-n sampling:

- randomly sample from P_t -
 → on each step t , restricted to just the top- n most probable words.
 → increasing n to get more diverse/risky output
 → Decrease n to get more generic/safe output
~~effie both of them are more efficient than beam search~~

Softmax temperature:

On timestep t , the LM computes a prob dist P_t by applying the softmax function to a vector of scores $s \in \mathbb{R}^M$

$$P_t(w) = \frac{\exp(s_w)}{\sum_{w' \in V} \exp(s_{w'})}$$

We can apply a temperature hyperparameter τ to the softmax

$$P_t(w) = \frac{\exp(s_w/\tau)}{\sum_{w' \in V} \exp(s_{w'}/\tau)}$$

→ Raise the temperature T : P_f becomes more uniform
Thus more diverse output (probability is spread around vocab)

→ Lower the temperature T : P_f becomes more spiky
Thus less diverse output

Softmax temperature is not a decoding algo
we can apply this during test time, in
conjunction with a decoding algorithm (such
as beam search or sampling)

Section 2: NLP tasks and neural approaches for them:

Summarization

Biggaword \rightarrow news article \rightarrow headline

LSTMs \rightarrow paragraph \rightarrow sentence summary

NYT, CNN/Daily Mail \rightarrow news article \rightarrow (multi) sentence summary

With how much full text? \rightarrow summary sentences

Sent:

- Sentence simplification is a different but related task
 - rewrite the source text in a simpler way
- simple Wikipedia: standard Wikipedia sentence → simple version
- Newsela: news article → version for children

Two main strategies:

Extractive summarization

Select parts (typically sentences) of the original text to form a summary

• Easier

• Restrictive (no paraphrasing)

Abstractive Summarization

Generate new text using natural language techniques

• More difficult

• More flexible (more human)

Evaluation:

ROUGE (Recall)-Oriented Understudy for Existing Evaluation

Evaluation

Rouge-N

like BLEU, it's based on n-gram overlap.

• ROUGE has no brevity penalty

• ROUGE based on recall, while BLEU is based on precision

• Precision is more important for MT and recall

is more important for summarization (assuming you have a max length constraint)

• However, often a F1 version of ROUGE is reported (anyway)

• BLEU is reported as a single number, which is combination of the precisions for $n=1, 2, 3, 4$ n-grams

• ROUGE scores are reported separately for each n-grams

ROUGE-L \Rightarrow unigram overlap

ROUGE-D \Rightarrow bigram overlap

ROUGE-LCS \Rightarrow LCS overlap

Neural approaches for summarization (2015-present)

→ 2015: Rush et al. publish the first seq2seq

Summarization paper

→ Single-document abstractive summarization is

a translation task!

Thus we can apply standard seq2seq + attention NMT methods

since 2015, developments

- Making it easier to copy [prevent too much of it]
- Hierarchical/multi-level attention
- More global/high-level content selection
- Using RL to directly maximize ROUGE on discrete goals
- Resurrecting pre-neural ideas (graph algos) and

working them ~~out~~ into neural systems

Copy Mechanism:

- Seq2Seq + attention are good at writing fluent output but bad at copying over details correctly
- Copy mechanisms use attention to enable seq2seq system to easily copy words and phrases from the input to the output.
 - Allowing both copying and generating gives us hybrid extractive/abstactive approach

$$P(w) = P_{\text{gen}} \text{Prob}(a) + (1 - P_{\text{gen}}) \sum_{i: w_i = w} \alpha_i^t$$

- Big problems with copying mechanism

They copy too much

bad at overall content selection

Preneural summarization had separate stages for
content selection and surface realization

Bottom-up summarization:

- Content selection stage: Use a neural sequence-tagging model to tag words as include or don't include
 - Bottom-up attention stage: The seq2seq+attention system can't attend to words tagged don't-include (apply a mask)
- Simple but effective → better overall content selection
Less copying

Neural summarization via RL:

- Use RL to directly optimize ROUGE-L
- ML+RL gives best result

Dialogue

• Task-oriented dialogue

→ Assistive (customer service, recommendations)

→ Co-operative (two agents solve a task together)

→ Adversarial (two agents compete in a task through dialogue)

• Social dialogue

→ chit-chat dialogue

→ Therapy/mental wellbeing

Preneural dialogue systems used predefined

templates or retrieve an appropriate response from a corpus of responses.

After 2015: open ended freeform dialogue.

→ Seq2Seq-based dialogue

However it quickly became apparent that a naive application of seq2seq has serious deficiencies for (chit chat) dialogue

- Genericness/boring responses
- Irrelevant responses (not related to experience)
- Repetition (showing the same thing over and over again)
- Lack of context
- Lack of consistent persona

Irrelevant response problem:

Reason:

- because it's generic
- or because changing the subject to something unrelated

One solution: optimize for MMI (Maximum Mutual Information)

between input S and response T

$$\log \frac{p(s, T)}{p(s)p(T)}$$

$$\hat{T} = \arg \max \left\{ \log \frac{p(T|s)}{p(T)} \right\}$$

Genericness or boring response problem

→ Easy test-time fixes:

- Directly upweight rare words during beam search
- Use a sampling decoding algo rather than beam search

→ Conditioning fixes

- Condition the decoder on some additional content (e.g. sample some content words and attend them to them)
- Train a retrieve-and-refine model rather than a generate-from-scratch model

Repetition Problems

Simple Solution

→ Directly block repeating n-grams during beam search

More complex solutions

- Train a coverage mechanism - in seq2seq
this is an objective that prevents the attention mechanism from attending to the same words multiple times.

- Define a training objective to discourage repetition
- If this is a differentiable function of the generated output, then will need some techniques e.g. RL to train

NLG: Storytelling:

Most neural storytelling work uses some kind of prompt:

- Generate a story-like paragraph given an image
- Generate a story given a brief writing prompt
- Generate the next sentence of a story given the story so far (story continuation)

Question: How to get around the lack of parallel data?

Answer: Use a common sentence encoding space

→ skip-thought vectors are a type of general-purpose sentence embedding method

- The idea is similar to how we learn an embedding for a word by trying to predict the words around it.

• Using COCO (an image captioning dataset), learn a mapping from images to the skip-thought encodings of their captions

• Using the target style corpus, train an RNN-LM to decode a skip-thought vector to the original text

• Put the two together

Generating story from a writing prompt

In 2018, Fan et. al. released a new story generation dataset collected from Reddit's WritingPrompts subreddit.

- Each story has an associated brief writing prompt or title.
 - propose a convolution-based story model
- Gated Multi-head multiscale self attention
- The self-attention is important for capturing long-range context
 - The gates allow the attention mechanisms to be more selective
 - The different attention heads attend to at different scales - this means there are different attention mechanisms dedicated to retrieving fine-grained information and coarse-grained information.

Model-fusion:

- Pretrain one seq2seq model, then train a second seq2seq model that has access to the hidden states of the first, with LM
- The idea is that the first seq2seq model learns general LM and the second learns to condition on the prompt.
- The results are impressive
 - Related to prompt
 - Diverse; non-generic
 - Stylistically dramatic
- However,
 - Mostly atmospheric/descriptive/scene-setting; less events/plot;
 - When generating for longer, mostly stays on the same idea without moving forward to new ideas - coherence issues.

NLG Evaluation:

Word overlap based metrics (BLEU, ROUGE, METEOR, F1, etc.)

- We know they are not ideal for machine translation worse for summarization and even worse for dialogue and storytelling

What about perplexity?

- Captures how powerful your LM is, but doesn't tell you anything about generation (e.g. if your decoding algo is bad, perplexity is unaffected)

Word embedding based metrics?

- compare the similarity of the word embeddings not just the overlap of the words themselves.
- Captures semantics in a more flexible way
- still doesn't correlate well with human judgements for open-ended task

Automatic evaluations of NLG

- We can define more focused automatic capture particular aspects of generated text
- Fluency (compute probability wrt a well trained LM)
- Correct style prob wrt ~~to the~~ LM trained on target corpus
- Diversity (rare word usage, uniqueness of strings, n-grams)
- Relevance to input (semantic similarity measures)
- Simple things like lengths and repetitions
- Task specific metrics e.g. compression rate for summarization

Though these doesn't measure overall quality, they can help us track some important qualities that we care about.

Human evaluation

regarded as gold standard but slow

and expensive

does human evaluation solve all your problems? No!

Conducting human evaluation effectively is very difficult

Humans are inconsistent

can be illogical

lose concentration

misinterpret your question

can't always explain why they feel the way they do.

Adversarial discrimination

- Test whether the NLG system can fool a discriminator which is trained to distinguish human text from artificially generated text

Section 4:

Thoughts on NLG research, current trends and the future:

- Incorporating discrete latent variables into NLG
- Alternatives to strict left-to-right generation
- Alternatives to maximum likelihood training with teacher forcing
- NLG is the wildest part remaining in NN
- Neural NLG community is rapidly expanding
- Biggest roadblock is still evaluation

8 things on NLG

- 1) The more open-ended the task the harder everything becomes
- 2) Aiming for specific improvement can also be more manageable than aiming to improve overall generation quality
- 3) If you're using a LM for NLG, improving the LM (i.e. perplexity) will most likely improve generation quality.
- 4) Look at your outputs a lot
- 5) You need an automatic metric, even if it's imperfect (probably several)
- 6) If we do human eval, make the questions as focused as possible
- 7) Reproducibility is a huge problem in today's NLP+DL, and a "huge" problem in NLG
• publicly release all your generated outputs.

Working with NLG is very frustrating.
But also very funny.

Non-autoregressive generation for NMT

In 2018 Gauvret et al. published Non-autoregressive generation for NMT. It generates the translation left-to-right, with each word depending on the ones before it.

- It generates the translation in parallel.
- It has obvious efficiency advantages but it is also intriguing from NLP point of view
- Transformer based architecture and decoder can run in parallel at text time

Lecture 16: Coreference Resolution:

- Identify all mentions that refers to the same real world entity

- * NLP models cannot handle this properly

Applications

- Full text understanding

- Information extraction, question answering

- Summarization

- "He was born in 1961" (Who?)

- Machine translation

- languages have different features for gender, number, dropped pronouns, etc.

- Dialogue systems

Coreference Resolution

- 1) Detect the mentions (easy)

- 2) Cluster the mentions (hard)

Mention Detection:

- Pronouns (pos tagger)

- Named entities (NER system)

- Noun phrases (con parsing)

Mention Detection: Not so simple

- Marking all pronouns, named entities and NPs as mentions

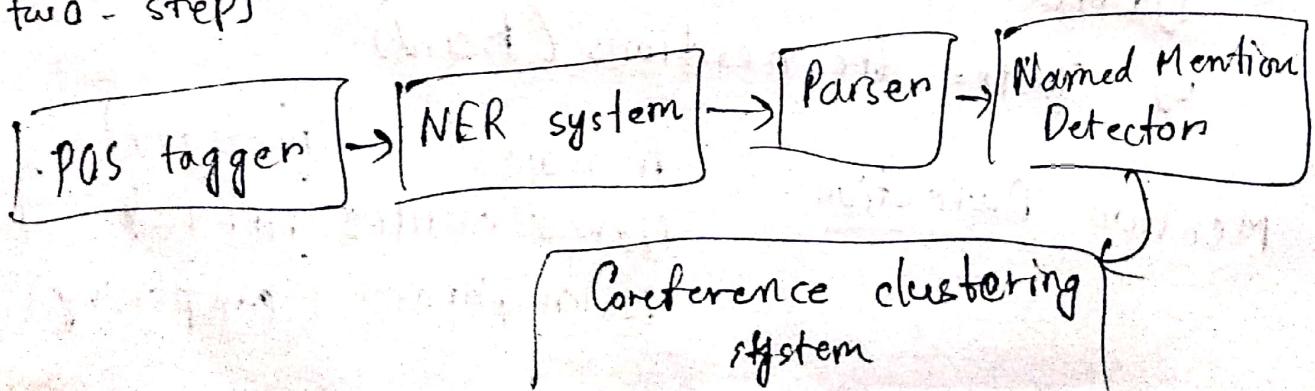
Are these mention?

- It is sunny
- Every student
- No student
- The best donut in the world
- 100 miles

How to deal with these bad mentions?

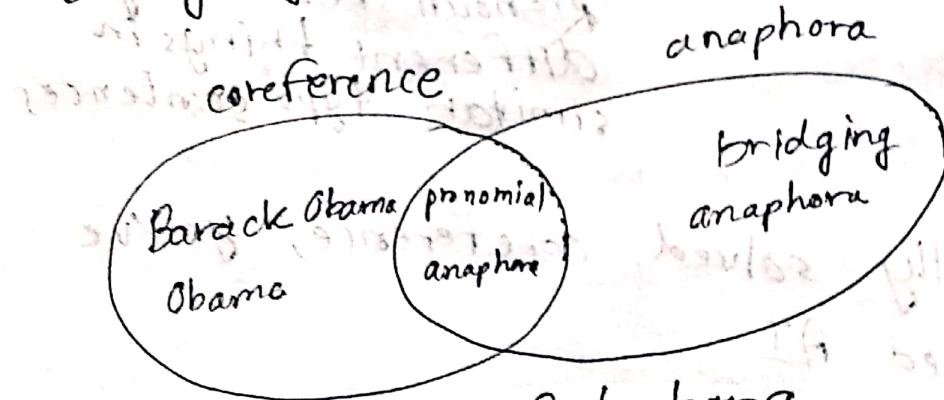
Can we avoid a pipelined system?

- We could instead train a classifier specifically for mention detection instead of using POS tagger, NER system, and parser.
- Or even jointly do mention-detection and coreference resolution end-to-end instead of two-steps



On to Coreference! First some linguistics

- Coreference is when two mentions refer to the same entity in the world!
- A related linguistic concept is anaphora: when a term (anaphor) refers to another term (antecedent)
 - The interpretation of the anaphor is in some way determined by the interpretation of the antecedent.
- Not all anaphoric relations are coreferential



Anaphora vs. Cataphora

- Usually the antecedent comes before the anaphor

Four Kinds of Coreference Models

- Rule-based (pronominal anaphora resolution).

- Mention Pair

Mention Ranking

- Clustering

- Traditional pronominal anaphora resolution:

Hobbs' naive algorithm (1976)

- Knowledge based Pronominal Coreference

Winograd (1972)

Winograd Schema

↓
pronoun referencing
different things in
similar type sentences

- If you've fully solved coreference, you're arguably solved AI

Mention Pair

Train a binary classifier that assigns every pair of mentions a probability of being coreferent

Mention Pair Models : Disadvantages

- Many mentions only have one clear antecedent
- But we are asking the model to predict all of them
- Instead train the model to predict only one antecedent for each mention
- More linguistically plausible

Coreference Models: Mention Ranking:

- Assign each mention its highest scoring candidate antecedent
→ Dummy NA mention allows model to decline linking the current mention to anything
- Positive examples: model has to assign a high probability to either one (but not necessarily both)
- Solve → Apply a softmax over the scores for candidate antecedents so probabilities sum to 1.

Training:

- We want the current mention m_j to be linked to any one of the candidate antecedents it's coreferent with.
- Mathematically, we want to maximize this probability

$$\sum_{j=1}^{i-1} \mathbb{1}(y_{ij}=1) p(m_j, m_i)$$

iterate through candidate antecedents (previously occurring mentions) for ones that are coreferent to m_j we want the model to assign high probability

- The model produces 0.9 probability for one of the correct antecedents and low probability for everything else, and the sum will be far still be large

$$J = \sum_{i=2}^N -\log \left(\sum_{j=1}^{i-1} \mathbb{1}(y_{ij}=1) p(m_j, m_i) \right)$$

usual trick of taking negative log to go from likelihood to loss.

Computing the probabilities

- a) Non-neural statistical classifier - (Extract features and use)
- b) Neural Coref Model (word embedding and few categorical features)
- c) End-to-end Model
SOTA [Lee, et.al. (W, EMNLP 2017)]
- LSTM
 - Attention
- for every subsequence, they come up with some span representation $g_i = [x_{\text{start}}^*, x_{\text{END}}^*, \hat{x}_i, \phi(i)]$

Attention scores $\alpha_t = w_a \cdot \text{FFNN}_a(x_t^*)$

$$s(i,j) = s_m(i) + s_m(j) + s_a(i,j)$$

↑ ↑ ↑
is that is that Do they look
a mention a mention coreferent

$$s_m(i) = w_m \cdot \text{FFNN}(g_i)$$

$$s_a(i,j) = w_a \cdot \text{FFNN}_a([g_i g_j, g_i \circ g_j, \phi(i, j)])$$

- Intractable to score
 - $O(T^2)$ spans of text in a document
 - $O(T^4)$ runtime

Clustering Based Approach

- Coreference is a clustering task. Let's use a clustering algorithm (agglomerative clustering) to merge each mention in its own singleton cluster.
- Start with each mention at each step.
- Merge a pair of clusters which cluster merges are good.
 - Use a metric to score
- Mention pair decision is difficult
- Cluster-pair decision is easier

MUC, CEAFF, LEA, B-CUBED, BLANC

Evaluation

operator

Conclusion:

- Coreference is a useful, challenging and linguistically interesting task
- Systems are getting better rapidly, largely due to better neural models

Lecture 17: Multi-task Learning
Richard Socher (chief scientist, Salesforce)

DecaNLP

Decathlon : Multi Learning
The Natural Language as Question Answering

Bryan McCann, Nitish Keshav, and
Joint work with

Caiming Xiong

[10 tasks in one model]

NLP and AI

Machine Learning
with feature
engineering

Deep Learning
for feature
engineering

Deep Architecture
engineering
for single tasks

what
next?

Single
Multitask
Model

Limits of Single Task learning:

- Great performance improvements in recent years given {dataset, task, model, metric}
- We can hill-climb to local optima as long as $|dataset| > 1000 \times C$
- For more of general AI, we need continuous learning in a single model instead
- Models typically start from random or are only partly pretrained

• Why has weight & model sharing not happened as much in NLP?

• Many types of reasoning:

→ NLP requires many types of reasoning: logical, linguistic, emotional, visual and others

logical, linguistic, emotional, visual and others

→ requires short and long term memory

→ requires short and long term memory

→ NLP had been divided into intermediate and separate tasks to make progress

• A single unsupervised task cannot solve it all

• Language clearly requires supervision in nature

• Language clearly requires supervision in nature

* Why a unified multi-task model for NLP?

* Why a unified multi-task model for NLP?

→ Multi-task learning is a blocker for general NLP systems

→ Unified models can decide how to transfer knowledge (domain adaption, weight sharing, transfer and zero shot learning)

→ Unified, multi-task models can

• more easily adapt to new tasks

• make deploying to production \times times simpler

• lower the bar for more people to solve new tasks

• Potentially move forward continual learning

How to express many NLP tasks in the same framework?

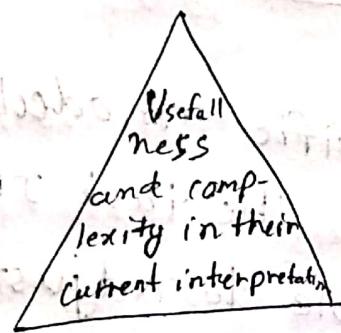
- Sequence tagging: NER, aspect specific sentiment
- Text classification: dialogue state tracking, sentiment classification
- Seq2Seq: machine translation, summarization, question answering

3 equivalent supertasks of NLP:

Language Modelling

Question Answering

Dialogue



- phrase extraction from context
- translation
- summary
- NLI
- sentiment classification
- semantic role labeling
- Relation extraction

• Dialogue state tracking

• SQL translation (semantic parsing)

• Winograd schema and anaphora resolution (common sense reasoning)

- Meta Supervised Learning: from $\{x, y\}$ to $\{x_t, y\}$ (t is the task)
- Use a question, q , as a natural description of the task, t , to allow the model to use linguistic information to connect tasks
- y is the answer to q and x is the context necessary to answer q

Designing a model for decaNLP

specifications

- No task-specific modules or parameters because we assume the task ID is not available.
 - Must be able to adjust internally to perform disparate tasks.
 - Should leave open the possibilities of zero-shot inference for unseen tasks.
- * Start with a context
- * Ask a question
- * Generate the answer one word at a time by
- Pointing to context
 - Pointing to question
 - Or choosing a word from an external vocabulary
- * Pointer switch chooses ~~between~~ among those three options

Multitask Question Answering Network (MQAN)

Initial Encoding → Co-attention → Transformer Layer × 2 →

Final Encoding

Datasets	Metrics	Task
SQuAD	nF1	QA
IWSLT En-De	BLEU	MT
CNN/Daily Mail	Rouge	Summarization
MultiNLI	EM	Natural Language Inference
SST-2	Sentiment Analysis	Sentiment
QA-SRT	nF1	Semantic Role Labeling
QA-ZRF	CF1	Relation Extraction
WOT	dS EM	Goal-Oriented Dialogue
WikiSQL	CFEM	Semantic Parsing
Winograd Schemas	EM	Pronoun Resolution

decScore = sum of all metrics

- Transformer layers yield benefit in single-task and multi-task learning

- QA and SRL have strong connection
- Pointing to the Question is essential
- Multitasking helps zero-shot

• There is a gap between the combined single-task models and the sig single multitask model.

④ Neural architecture search
→ Reinforcement learning to figure out what layer comes next

Training strategy

Fully joint \rightarrow Round-Robin build minibatch out of all tasks.

Anti-curriculum pretraining

→ train on the harder tasks and add the simpler tasks later on

Closing the Gap

i) Fully joint training

ii) Anti-curriculum pretraining

iii) GloVe to CoVe

iv) including more tasks: anti-curriculum pretraining

v) oversampling IWSLT

When MQAN points

- Answers are correctly copied from either context or question
- No confusion over which task the model should perform on which output space to use.

Zero shot Domain Adaptation on pretrained MQAN

- achieves 80% on Amazon and Yelp reviews

, achieving 80% on Amazon and Yelp reviews

- Achieves 62% on SNLI (87% with fine-tuning, a

2-point gain over random initialization)

- BERT is not the silver-bullet to multi-task learning

, BERT is not the silver-bullet to multi-task learning

Most exciting thing would be combining multiple tasks.

- If our model can translate English to German

and summarize, why can't we want a

German summary? (interesting question)

decaNLP is a benchmark for generalized NLP

- Train single question answering model for multiple NLP tasks

Framework for tackling NLP tasks

→ more general language understanding

→ multi-task learning

→ domain adaptation

→ Transfer learning

→ weight sharing, pre-training, fine-tuning

→ zero-shot learning