

CMPT 413/713: Natural Language Processing

### Project tips and analyzing your results

Spring 2025 2025-03-12

### Project Milestone

- Project Milestone due Thursday 3/20
- PDF (3-6 pages) in the style of a conference (e.g. ACL/EMNLP) submission • https://2020.emnlp.org/files/emnlp2020-templates.zip
- Milestone should include:
  - **Title and Abstract** summary of what you are working on

  - **Introduction** motivate the problem, describe your goals, and highlight your findings • **Prior Work** - what have others done in this area?
  - Approach details on your main approach and baselines. Be specific. Make clear what part is original, what code you are writing yourself, what code you are using
  - **Experiment** describe dataset, evaluation metrics, what experiments you plan to run, any results you have so far. Also provide training details, training times, etc.
  - Future Work what is your plan for the rest of the project
  - **Reference** provide references using BibTex
- Milestone will be graded based on **progress** and **writing quality**

### Final project report

Rough page layout: your report can have different number of pages for each section

Title Abstract Introduction	Prior related work
Data	Experiment

8 pages not including references

Model or what approach you are taking



# Me expect you to have around 3-6 pages, can be more (but < 8) if you

We expect you to have around 3-6 pages, can be more (but < 8) if you have done lots of work since the proposal



Have statistics/ analysis of your data, and show examples!

### Include References

- Good writing is concise and clear
- More words does not mean more information or higher quality
- Trim words you don't need.
- But I have a lot of stuff to say!
  - References don't count toward the page limit
  - appendices will not be graded)

### Why a page limit?

• You can have appendices if you really want to share (the

### Project Milestone

- Build on what you have done for the proposal
- Include progress and initial results
- Flesh out description of your approach
- Include figures for your model
- Include dataset examples and statistics
- Have plan for experiments and analysis you will do
- Include references
- Make it clear what you are implementing vs what part you are building on top of existing libraries / codebases / homework

### What you should have from the Project Proposal:

- What **task** are you addressing? What is the **input / output**? Why is it interesting?
- What specific aspects will your project be on?
  - Re-implement paper? Compare different methods? Analysis?
- What have others (**prior work**) done to address the same problem?
- What **data** do you plan to use?
  - Preliminary statistics for your data (number of sentences, tokens, etc)
- What is the specific method or methods you will use to address the task?
  - What will you implement by yourself vs what existing code will you use?
  - What **compute resources** do you plan to use?
- How do you plan to evaluate?
  - Data splits?
  - What **metrics**?
  - What experiments will you run to **compare** different variations / different approaches?
- Timeline and work breakdown
  - What do you plan to have by the milestone? The end of the term?
  - Who will work on what?

### Project Milestone

Build on what you have done for the proposal

- Task / Problem State:
  - Clearly state the input and output
- **Related Work:** How have others approached this problem?
- - existing libraries / codebases / homework
  - use?
  - Make figures (if appropriate) that illustrate your approach.

• **Progress**: What have you achieved so far? Are there any issues you encountered?

• Approach: Flesh out and update your approach based on what you have learned.

• Make it clear what you are implementing vs what part you are building on top of

• Specify what you have implemented and what you still have to implement.

• Focus on describing how you are using a particular model for your task. For instance, if you are using a RNN for text classification, explain how the input is fed into the RNN, and how the RNN is used to make predictions. What is the training loss that you will

• Using equations (if appropriate) with clear mathematical style to explain your model(s).





# Project Milestone

Build on what you have done for the proposal

- **Data:** At this point, you should know what data you will be using. Describe the data and provide some statistics and examples from the data.
- Experiments and results: At this point, you should have a clear idea of what experiments you will perform and/or some preliminary results.
  - Provide a list of experiments you have will perform. Describe what you expect the experiments to reveal, or what is uncertain about the potential outcomes.
  - Provide a summary of your preliminary results, and describe remaining results that you plan to produce
  - Have a empty table with rows and columns
- **Timeline / Plans for remaining tasks**: Present an updated timeline of the planned tasks/goals.
  - Clearly state what you plan to complete by the final report.
- If you are working in a group, please also state the contribution of each team member • **References:** Many of you already had references in the proposal. For the milestone, this is required.



### • Task

- that
- Simple and effective way to do it is via an example
- Summarize **relevant work** 
  - Summarize what did they do? What was the key findings?

  - plagiarism).
  - structures.

### Proposal Observations

• There are standard NLP tasks: sentiment analysis, question answering • Tasks are defined by their input and output - so make sure you clearly specify

• (this requires larger changes than just replacing "we" with "they") • Avoid copying text from another paper verbatim (if you copy, it is

• Describe the approach in your own words. Use different sentence

### Task / method illustration example

### Clear input

Q: How many of Warsaw's inhabitants spoke Polish in 1933?



WikipediA

The Free Encyclopedia

### Document Retriever



100.14

No. 10, No. 104 (Sec.)

#### Warson

Page Miligada, He has propriogenia

No cettin è giografite fisisi aggini l'evoler con con Nover (donnigati

"Bennet' substitues for the size on Branes (herity-and) "Do rivieser' estes for i'vie louret iter We been paster, print, 34 Ante Spin-Spanter, 3 The an interact builds

NAME AND INCOMES OF ANY ADDRESS OF ADDRESS OF ADDRESS OF ADDRESS ADDRES ADDRESS stands on the Wester New Accession (Allera) angles (Millionation (Wiles) And the balls frequent UK INCOMENTATION AND A DESCRIPTION OF A DES contraction must be interested, which eater theme in the set parties upto the later Answer war war was not the second of the second of the second sec come 4, 10, 40 ages and sharehow 1, 200, 20 ages (17

1.418 for her south of the spectra of the spectra disk had been handle into a fact and <sup>24</sup> have also control or control from many interesting into the control of program. Topics, Manuae in control or a first or a first of a control of the con have the second to be had and a shore place where without a second ball." "" these has a to a nite state of relation in the stretter of 1600 monifold in a minimum relation of rendering at his pressing. The dynamics is supported to the set of an analysis with the second as your beau really appearing the second balance of the bages of man supreme to be an en-Index decays "Vision Reducings The agency is adored units and prior in temporary of NAME AND DESCRIPTIONS, SAME OF PARTY AND ADDRESS AND ADDRESS ADDRE The highest section of appropriate to be harped without "" Minister in other two wind Vision documents." LANS LINE OF THE ALL AND REAL AND REA



#### Reading Wikipedia to Answer Open-Domain Questions [Chen et al, 2017]



### • Approach / Method / Model

- This is what you proposed to do to solve the task
  - Describe how you will use a particular method for your task
  - It's great that you are using RNN, LSTM, GRU, and that transformers uses self-attention, and know the differences between the models.
  - But how will you actually use it for your task?
- Typically in NLP we use **neural network models**, so often it is
  - a description of the model architecture with information on
  - how the **input** is fed in
  - how the **prediction** is made
  - how the model is trained (**optimizer and loss function**)

### Proposal Observations

## Method figure example

# Figure showing key elements

#### (a) Unsupervised SimCSE

Different hidden dropout masks in two forward passes



Figure 1: (a) Unsupervised SimCSE predicts the input sentence itself from in-batch negatives, with different hidden dropout masks applied. (b) Supervised SimCSE leverages the NLI datasets and takes the entailment (premise-hypothesis) pairs as positives, and contradiction pairs as well as other in-batch instances as negatives.

Training objective (loss function)

$$\begin{split} \ell_i = -\log \frac{e^{\min(\mathbf{h}_i^{z_i}, \mathbf{h}_i^{z_i'})/\tau}}{\sum_{j=1}^N e^{\sin(\mathbf{h}_i^{z_i}, \mathbf{h}_j^{z_j'})/\tau}} \end{split}$$

SimCSE: Simple Contrastive Learning of Sentence Embeddings [Gao et al, 2022]

#### (b) Supervised SimCSE

### • If you are doing **prompt engineering** with large language models

- Clearly describe the prompts you are trying out
- If you are giving the LLM examples, make sure to clearly specify how the examples are specified
- Clearly specify how you will interpret the LLM output.
  - Do you see errors from the LLM?
  - How do you handle unexpected responses from the LLM
- Be systematic

### Proposal Observations

### Clear input

#### Standard Prompting

#### Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Clear output

Model Output

A: The answer is 27. 🗙

Figure 1: Chain-of-thought prompting enables large language models to tackle complex arithmetic, commonsense, and symbolic reasoning tasks. Chain-of-thought reasoning processes are highlighted.

Chain-of-Thought Prompting Elicits Reasoning in Large Language Models [Wei et al, 2022]



### Clear difference in prompting





- Experiments
  - The quality of a model / approach is measured by **evaluation metrics**
  - For machine learning models, you will train and evaluate them with **data**. Specify clearly your data and statistics about your data.
  - Typically you want have a set of **comparisons** (often informed by some hypothesis) • There can be **variations** on your model depending on hyper-parameters, input encoding, how it is trained (training data, optimizer, loss function), etc. • Hypothesis: Adam converges faster than SGD. -> Experiment: Train with Adam and
    - - SGD and compare.
      - models for text classification -> Experiment: Train Bi/Uni-dir models and compare.
    - Hypothesis: Bidirectional RNNs gives higher performance than unidirectional • You can also compare different models
  - Try to be **concrete** about your plans

### Proposal Observations

- Use the correct format
- References
  - Use bibtex and \cite commands
  - (word2vec, BERT, LSTM, GRU, Transformer, etc).
  - Make sure that the references are properly included
- English
  - Use complete sentences
  - Proofread your report
  - Have friend (native/fluent English speaker) proofread your report

### Proposal Observations

• Make sure that you cite the papers that introduced the method you are using

### Timeline

- your model.
- classifier, and analyzing the results.
- Plan for report writing and video making

• Make sure to allocate time for training and debugging the training of

• In addition, you should allocate time for setting up the experiments (train classifier with some set of hyperparameters), evaluate the

# Tips for good final projects

• Have a clear, well-defined hypothesis to be tested

(++ novel/creative hypothesis)

- Conclusions and results should teach the reader something
- Meaningful tables, plots to display the key results

++ nice visualizations or interactive demos

++ novel/impressive engineering feat

++ good results

### What to avoid

- All experiments run with prepackaged source no extra code written for model/data processing
- Just ran model once or twice on the data and reported results (not much hyperparameter search done)
- A few standard graphs: loss curves, accuracy, without any analysis
- Results/Conclusion don't say much besides that it didn't work
  - Even if results are negative, you should analyze them!

Remember: Include analysis!

# What makes for a good paper? (example paper)

A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task <u>https://aclanthology.org/P16-1223.pdf</u> (Chen et al, ACL 2016, Outstanding Paper)

# Task definition

### Example with input and output

#### Text describing the problem

an example<sup>4</sup>: it consists of a passage p, a question q and an answer a, where the passage is a news article, the question is a cloze-style task, in which one of the article's bullet points has had one entity replaced by a placeholder, and the answer is this questioned entity. The goal is to infer the missing entity (answer a) from all the possible entities which appear in the passage. A news article is usually

(@entity4) if you feel a ripple in the force today, it may be the news that the official @entity6 is getting its first gay character . according to the sci-fi website @entity9, the upcoming novel " @entity11 " will feature a capable but flawed @entity13 official named @entity14 who " also happens to be a lesbian . " the character is the first gay figure in the official @entity6 -- the movies, television shows, comics and books approved by @entity6 franchise owner @entity22 -- according to @entity24 , editor of " @entity6 " books at @entity28 imprint @entity26 .

#### Question

characters in " @placeholder " movies have gradually become more diverse

A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task https://aclanthology.org/P16-1223.pdf (Chen et al, ACL 2016, Outstanding Paper)

#### Passage

Taken from the Reading Comprehension dataset introduced by Herman et al, 2015

### Answer

@entity6

Goal: identify entity that is goes where @placeholder goes



# Task definition

- Dataset is automatically generated
  - Use Google NLP pipeline and get entities and coreference chains
  - Cloze-style (fill in the blank)  $\bullet$ questions generated by taking sentence from passage and replacing entity reference with @placeholder
  - Data is anonymized  $\bullet$ (entities are just @entity#)

### Example with input and output

(@entity4) if you feel a ripple in the force today, it may be the news that the official @entity6 is getting its first gay character . according to the sci-fi website @entity9 , the upcoming novel " @entity11 " will feature a capable but flawed @entity13 official named @entity14 who " also happens to be a lesbian . " the character is the first gay figure in the official @entity6 -- the movies, television shows, comics and books approved by @entity6 franchise owner @entity22 -- according to @entity24 , editor of " @entity6 " books at @entity28 imprint @entity26 .

#### Question

characters in " @placeholder " movies have gradually become more diverse

A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task https://aclanthology.org/P16-1223.pdf (Chen et al, ACL 2016, Outstanding Paper)

#### Passage

Taken from the Reading Comprehension dataset introduced by Herman et al, 2015



## Anonymization

#### **Original Version**

#### Context

The BBC producer allegedly struck by Jeremy Clarkson will not press charges against the "To Gear" host, his lawyer said Friday. Clarkson, w hosted one of the most-watched television show in the world, was dropped by the BBC Wedness after an internal investigation by the British brocaster found he had subjected producer Oisin " "to an unprovoked physical and verbal attack."

#### Query

Producer X will not press charges against Jeremy producer X will not press charges against *ent212*, Clarkson, his lawyer says. his lawyer says.

#### Answer

Oisin Tymon

Teaching machines to read and comprehend <u>https://arxiv.org/pdf/1506.03340.pdf</u> (Hermann et al, NIPS 2015)

#### **Anonymised Version**

	(1 (207 1 1) 1) (11 (270 1)
1	the <i>ent381</i> producer allegedly struck by <i>ent212</i> will
op	not press charges against the "ent153" host, his
who	lawyer said friday . ent212, who hosted one of the
WS	most - watched television shows in the world, was
sday	dropped by the ent381 wednesday after an internal
oad-	investigation by the ent180 broadcaster found he
Tymon	had subjected producer ent193 "to an unprovoked
,	physical and verbal attack . "

ent193



### Two approaches presented

In this section, we describe two systems we implemented – a conventional entity-centric classifier and an end-to-end neural network. While Hermann

Features are described for entitycentric classifier

> A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task https://aclanthology.org/P16-1223.pdf (Chen et al, ACL 2016, Outstanding Paper)

### Approach

Describe problem again (with math symbols)

Given the (passage, question, answer) triple  $(p, q, a), p = \{p_1, \dots, p_m\}$  and  $q = \{q_1, \dots, q_l\}$  are sequences of tokens for the passage and question sentence, with q containing exactly one "@placeholder" token. The goal is to infer the correct entity  $a \in p \cap E$  that the placeholder corresponds to, where E is the set of all abstract entity markers. Note that the correct answer entity must appear in the passage p.

# Model description

- Describe key parts including how it is hooked out to input and how prediction is made
- **Encoding:** First, all the words are mapped to *d*-dimensional vectors via an embedding matrix  $E \in \mathbb{R}^{d \times |\mathcal{V}|}$ ; therefore we have *p*:  $\mathbf{p}_1, \ldots, \mathbf{p}_m \in \mathbb{R}^d$  and  $q : \mathbf{q}_1, \ldots, \mathbf{q}_l \in \mathbb{R}^d$ .

**Prediction:** Using the *output* vector **o**, the system outputs the most likely answer using:

$$a = \arg \max_{a \in p \cap E} W_a^{\mathsf{T}} \mathbf{0} \tag{4}$$

Symbols are connected

A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task <u>https://aclanthology.org/P16-1223.pdf</u> (Chen et al, ACL 2016, Outstanding Paper)

# Model architecture with input and output

Passage

(@entity4) if you feel a ripple in the force today, it may be the

news that the official @entity6 is getting its first gay character . according to the sci-fi website @entity9 , the upcoming novel " @entity11 " will feature a capable but flawed @entity13 official

named @entity14 who " also happens to be a lesbian . " the character is the first gay figure in the official @entity6 -- the

movies , television shows , comics and books approved by @entity6 franchise owner @entity22 - according to @entity24 , editor of " @entity6 " books at @entity28 imprint @entity26 .

#### Question





## Provide implementation and training details

- What external software / code was used
- Why was it used?

### Training Details (for entity-centric model)

For training our conventional classifier, we use the implementation of *LambdaMART* (Wu et al., 2010) in the RankLib package.<sup>5</sup> We use this ranking algorithm since our problem is naturally a ranking problem and forests of boosted decision trees have been very successful lately (as seen, e.g., in many recent Kaggle competitions). We do not use all the features of *LambdaMART* since we are only scoring 1/0 loss on the first ranked proposal, rather than using an IR-style metric to score ranked results. We use Stanford's neural network dependency parser (Chen and Manning, 2014) to parse all our document and question text, and all other features can be extracted without additional tools.

## Provide implementation and training details

### Training Details (for NN model)

- Vocabulary size, UNK handling
- Word embeddings
- Parameter initialization
- Dimensions (LSTM hidden state, word embeddings)
- Optimizer: type, learning rate, minibatch size, dropout
- Number of epochs trained, how the "best" model was selected
- Compute resource uses (GPU type, runtime)

vanilla stochastic gradient descent (SGD), with a For training our neural networks, we only keep fixed learning rate of 0.1. We sort all the examples the most frequent  $|\mathcal{V}| = 50k$  words (including enby the length of its passage, and randomly sample a tity and placeholder markers), and map all other mini-batch of size 32 for each update. We also apply words to an *<unk>* token. We choose word embeddropout with probability 0.2 to the embedding layer ding size d = 100, and use the 100-dimensional preand gradient clipping when the norm of gradients trained GloVe word embeddings (Pennington et al., exceeds 10. 2014) for initialization. The attention and output pa-All of our models are run on a single GPU rameters are initialized from a uniform distribution (GeForce GTX TITAN X), with roughly a runtime between (-0.01, 0.01), and the LSTM weights are of 6 hours per epoch for CNN, and 15 hours per initialized from a Gaussian distribution  $\mathcal{N}(0, 0.1)$ . epoch for *Daily Mail*. We run all the models up to 30 epochs and select the model that achieves the We use hidden size h = 128 for CNN and 256 best accuracy on the development set. for Daily Mail. Optimization is carried out using





### Provide implementation and training details

	Hyperparam	BERT	RoBERTa	ALBERT
	Epochs	3,10,20	3	3
	Learning rate	1e - 5 - 5e - 5	1e - 5 - 3e - 5	1e - 5 - 3e - 5
	Learning rate schedule	warmup-linear	warmup-linear	warmup-linear
Organize hyper	Warmup ratio	0.1	0.1	0.1
narameters in a table	Batch size	16	16	16
	Adam $\epsilon$	1e-6	1e-6	1e-6
if appropriate	Adam $\beta_1$	0.9	0.9	0.9
	Adam $\beta_2$	0.999	0.98	0.999
	Adam bias correction	{True, False}	{True, False}	{True, False}
	Dropout	Ò.1	0.1	_
	Weight decay	0.01	0.1	_
	Clipping gradient norm	1.0	_	1.0
	Number of random seeds	25	25	25

On the stability of fine-tuning BERT: Misconceptions, explanations, and strong baselines <u>https://arxiv.org/pdf/2006.04884v3.pdf</u> (Mosbach et al, ICLR 2021)

### Provide data statistics

	379,450
# Train 380,298	
<ul> <li>Statistics about what goes</li> <li>into train/val/test</li> <li># Dev</li> <li>3,924</li> </ul>	64,835
# Test 3,198	53,182
• NLP specific statistics Passage: avg. tokens 761.8	813.1
(number of tokens, Passage: avg. sentences 32.3	28.9
sentences, unique words, Question: avg. tokens 12.5	14.3
Avg. # entities 26.2	26.2

A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task <u>https://aclanthology.org/P16-1223.pdf</u> (Chen et al, ACL 2016, Outstanding Paper)

### Summarize your results in tables

Test 40.2 50.9 57.0 63.0	Dev 35.5 56.4 63.3 70.5	Test 35.5 55.5 62.2 69.0
40.2 50.9 57.0 63.0	35.5 56.4 63.3 70.5	35.5 55.5 62.2 69.0
50.9 57.0 63.0	56.4 63.3 70.5	55.5 62.2 69.0
57.0 63.0	63.3 70.5	62.2 69.0
63.0	70.5	69.0
62.8	40.0	
05.0	69.0	68.0
60.6	N/A	N/A
66.8	N/A	N/A
69.4*	N/A	N/A
67.9	69.1	68.3
72.4	76.9	75.8
	66.8 69.4* 67.9 <b>72.4</b>	66.8       N/A         69.4*       N/A         67.9       69.1         72.4       76.9

- Use bilinear attention instead of tanh layer
- Simpler model (only use weights contextual embeddings, don't classify over all vocabulary but only words that appear in passage)

ININ MODELIS SIMILAR TO ALLEMINE Reader, Compared to Allemine Reader

# Conduct ablation studies

to ablate = to remove

- If classical classifier (Naive Bayes, logistic regression, decision forest), do feature ablation to see how important is each feature
- If neural network with different components, ablate components.
- If loss is combination of terms, ablate loss terms.

A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task <u>https://aclanthology.org/P16-1223.pdf</u> (Chen et al, ACL 2016, Outstanding Paper)



Features	Accuracy
Full model	67.1
- whether $e$ is in the passage	67.1
- whether $e$ is in the question	67.0
- frequency of $e$	<b>63.7</b>
<ul> <li>position of e</li> </ul>	65.9
<ul> <li><i>n</i>-gram match</li> </ul>	60.5
<ul> <li>word distance</li> </ul>	65.4
<ul> <li>sentence co-occurrence</li> </ul>	66.0
<ul> <li>dependency parse match</li> </ul>	65.6

### Perform data analysis

- Manual analysis of 100 samples from val (dev) split
- Categorize into 6 different types (based on relationship of passage, question and answer)

No.	Category	(%)
1	Exact match	13
2	Paraphrasing	41
3	Partial clue	19
4	Multiple sentences	2
5	Coreference errors	8
6	Ambiguous / hard	17

A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task https://aclanthology.org/P16-1223.pdf (Chen et al, ACL 2016, Outstanding Paper)

- Answerable categories:
  - Exact match
  - Paragraphing
  - Partial clue
  - Multiple sentences
- Unanswerable categories:
  - Coreference error
    - The dataset was automatically generated
  - Ambiguous or very hard
    - Humans are not likely to be able to answer these



# Concrete examples for each category

Category	Question	]
Exact Match	<i>it 's clear @entity0 is leaning to- ward @<b>placeholder</b>, says an ex- pert who monitors @entity0</i>	
		ä
Para- phrase	@placeholder says he under- stands why @entity0 wo n't play at his tournament	]
Partial clue	a tv movie based on @entity2 's book @ <b>placeholder</b> casts a @en- tity76 actor as @entity5	i
Multiple sent.	he 's doing a his - and - her duet all by himself, @entity6 said of @placeholder	, 1
Coref. Error	rapper @ <b>placeholder</b> " disgusted , " cancels upcoming show for @en- tity280	1 1
Hard	pilot error and snow were reasons stated for <b>@placeholder</b> plane crash	( 1

#### Passage

...@entity116, who follows @entity0's operations and propaganda closely, recently told @entity3, *it 's clear @entity0 is leaning toward @entity60* in terms of doctrine, ideology and an emphasis on holding territory after operations....

... @entity0 called me personally to let me know that he would n't be playing here at @entity23, " @entity3 said on his @entity21 event 's website ....

... to @entity12 @entity2 professed that his @entity11 is not a religious book ....

... we got some groundbreaking performances, here too, tonight, @entity6 said. we got @entity17, who will be doing some musical performances. he's doing a his and - her duet all by himself....

... with hip - hop star @entity246 saying on @entity247 that he was canceling an upcoming show for the @entity249 . ... (but @entity249 = @entity280 = SAEs)

... a small aircraft carrying @entity5, @entity6 and @entity7 the @entity12 @entity3 crashed a few miles from @entity9, near @entity10, @entity11....

# Analysis leads to insights

	How many can the two models	(%)	Category	C	lassifier	Ne	ural net
	get correct for each category?	13	Exact match	13	(100.0%)	13	(100.0%)
		41	Paraphrasing	32	(78.1%)	39	(95.1%)
	Neural network is at close to	19	Partial clue	14	(73.7%)	17	(89.5%)
	optimal performance on this dataset!	2	Multiple sentences	1	(50.0%)	1	(50.0%)
		8	Coreference errors	4	(50.0%)	3	(37.5%)
		17	Ambiguous / hard	2	(11.8%)	1	(5.9%)
			All	66	(66.0%)	74	(74.0%)

### Analyzing your results

## Analyzing your results

- Understanding the output of your model
  - getting wrong
    - Provide examples
  - Conduct data / error analysis
  - Visualize
- Characterizing your model performance
  - Compare against other models
  - Compare against variations
    - Ablation studies

• Look at the output and examine what it is getting right and what it is

### Provide qualitative examples of your model output

• Simple tip: color code correct and incorrect output

#### **English-German translations**

src	Orlando Bloom and Miranda Kerr still love of
ref	Orlando Bloom und Miranda Kerr lieben sic
best	Orlando Bloom und Miranda Kerr lieben ein
base	Orlando Bloom und Lucas Miranda lieben
src	"We' re pleased the FAA recognizes that an
	with safety and security, " said Roger Dow
ref	"Wir freuen uns , dass die FAA erkennt , da
	spruch zur Sicherheit steht ", sagte Roger D
best	" Wir freuen uns , dass die FAA anerkennt
	Sicherheit unvereinbar ist ", sagte Roger Do
base	" Wir freuen uns über die <unk>, dass ein «</unk>
	Sicherheit und Sicherheit ", sagte Roger Ca

Effective Approaches to Attention-based Neural Machine Translation https://arxiv.org/pdf/1508.04025.pdf (Luong et al, 2015)

each other ch noch immer nander noch immer. einander noch immer. n enjoyable passenger experience is not incompatible

, CEO of the U.S. Travel Association . ss ein angenehmes Passagiererlebnis nicht im Wider-

ow, CEO der U.S. Travel Association.

, dass ein angenehmes ist nicht mit Sicherheit und ow, CEO der US - die.

<unk> <unk> mit Sicherheit nicht vereinbar ist mit meron, CEO der US - <unk>.

### Conduct data and error analysis

How to perform data and error analysis?

How to start?

- Take a manual subsample of the data • Manually group and categorize them (like the paper we
- looked at)
  - You will need to figure out what these groups are!
- Compute some statistics
  - Overall percentage of categories
  - What is the performance on each category?

Adapted from slide by Graham Neubig



### Conduct data and error analysis

How to perform data and error analysis?

General recipe

- Partition the performance of the validation set into attributes.
  - Define attributes
  - Group test samples
  - Breakdown of performance

different interpretable grouped based on pre-defined

Adapted from slide by Graham Neubig



### Defining attributes

- Different tasks could have different attributes
- Token-level, span-level, sentence-level
  - Token-level: part-of-speech tag
  - Span-level: span length
  - Sentence-level: sentence length

Slide credit: Graham Neubig



#### Example: breakdown performance by entity length Performance Histogram $\phi_{eLen} = 2$ F1... ... ... New eLen... York $\geq 4$ 3 2eLen... • • • $2 \quad 3 \geq 4$ 1 ... Bucketing Attributes Breakdown



Slide credit: Graham Neubig



• Compare across attributes



• Compare across models









#### Performance Gap Histogram





#### Slide credit: Graham Neubig



# Comparison of model performance by sentence length

• Attribute: sentence length



Effective Approaches to Attention-based Neural Machine Translation <u>https://arxiv.org/pdf/1508.04025.pdf</u> (Luong et al, 2015)

### Integrating unit tests

- Small careful test sets sound like... unit test suites, but for neural networks!
- Minimum functionality tests: small test sets that target a specific behavior.



- ٠ with categories of linguistic capabilities and types of tests.

Beyond Accuracy: Behavioral Testing of NLP models with CheckList https://arxiv.org/pdf/2005.04118.pdf (Ribeiro et al, ACL 2020)

	Expected	Predicted	Pass?			
<b>AFT</b> La	bels: negativ	ve, positive,	neutral			
{POS_VERB} the {THING}.						
e food.	neg	pos	X			
	neg	neutral	X			

Failure rate = 76.4%

Ribeiro et al., 2020 showed ML engineers working on a sentiment analysis product an interface

The engineers found a bunch of bugs (categories of high error) through this method!

Slide credit: John Hewitt



# Comparisons

# Guidelines for making comparisons

- Make sure you have baselines
  - Random
  - Majority Class
  - Conditioned Majority Class
- Other baselines
  - Simpler version of the model vs more complex
    - One layer FFNN
    - Unidirectional/Single-layer RNN
  - Retrieval baseline for generation problems

- Always consider: is the complexity warranted?
  - Different types of complexity
    - Simplicity in implementation vs simplicity in concept vs simplicity in model (number of parameters)
  - NN vs rule based vs NBs vs logistic regression

### Classifier comparison

- - Practical differences
  - Confidence intervals
  - Wilcoxon signed-rank test (covered in the sentiment unit)
  - McNemar's test (covered in the sentiment unit)

• Suppose you've assessed two classifier models. Their performance is probably different to some degree. What can be done to establish whether these models are different in any meaningful sense?

Slide credit: Chris Potts

### Multiple runs with different seeds



On the stability of fine-tuning BERT: Misconceptions, explanations, and strong baselines <u>https://arxiv.org/pdf/2006.04884v3.pdf</u> (Mosbach et al, ICLR 2021)

# Tips for training and debugging

- TA Tutorial
- Stanford CS231n:
  - <u>https://cs231n.github.io/neural-networks-2/</u>
  - https://cs231n.github.io/neural-networks-3/
- CMU nn4nlp: Graham Neubig

  - <u>https://www.youtube.com/watch?v=KRQHdwpfj-4</u>

http://www.phontron.com/class/nn4nlp2021/schedule/debugging.html

# Good luck on your project!