### Machine Translation and \* Advanced Deep Learning Models .\* • for Sequence Processing

Natalie Parde UIC CS 521 Review: Essentials of Deep Learning Feedforward neural networks

Convolutional neural networks

Recurrent neural networks

Stacked models

Bidirectional models







Natalie Parde - UIC CS 521



Natalie Parde - UIC CS 521





Natalie Parde - UIC CS 521



### Why does this work?

 When computational units are combined, the outputs from each successive layer provide new representations for the input that can better separate the data into the target classes



### Why does this work?

 When computational units are combined, the outputs from each successive layer provide new representations for the input that can better separate the data into the target classes



Review: Essentials of Deep Learning Feedforward neural networks

Convolutional neural networks

Recurrent neural networks

Stacked models

Bidirectional models

#### Convolutional Neural Networks

- Neural networks that incorporate one or more convolutional layers
- Designed to reflect the inner workings of the visual cortex system
- Require that fewer parameters are learned relative to feedforward networks for equivalent input data

In NLP, convolutions are typically performed on entire rows of an input matrix, where each row corresponds to a word.



Natalie Parde - UIC CS 521

In NLP, convolutions are typically performed on entire rows of an input matrix, where each row corresponds to a word.



Natalie Parde - UIC CS 521













### It's common to extract multiple different feature maps from the same input.



Natalie Parde - UIC CS 521

• • • • • • • • • • • •

### **Pooling Layers**



### The output from pooling layers is passed along as input to the rest of the network.



Review: Essentials of Deep Learning Feedforward neural networks

Convolutional neural networks

Recurrent neural networks

Stacked models

Bidirectional models

Recurrent Neural Networks (RNNs) Built-in capacity to handle temporal information

 Contain cycles within their connections, where the value of a unit is dependent upon outputs from previous timesteps

Can accept variable length inputs without the use of fixed-size windows

#### Many varieties exist

- "Vanilla" RNNs
- Long short-term memory networks (LSTMs)
- Gated recurrent units (GRUs)









Sequence processing models like **RNNs** are also useful for many classification problems.

- Sequence Labeling Tasks: Given a fixed set of labels, assign a label to each element of a sequence
  - Example: Part-of-speech tagging
- Inputs  $\rightarrow$  word embeddings
- Outputs → label probabilities generated by the softmax (or other activation) function over the set of all labels
- Sequence Classification Tasks: Given an input sequence, assign the entire sequence to a class (rather than the individual tokens within it)

### **Sequence Labeling**



### **Sequence Labeling**





### How to use RNNs for sequence classification?



### **Sequence Classification**



### **Sequence Classification**



Natalie Parde - UIC CS 521

#### Long Short-Term Memory Networks (LSTMs)

- Remove information no longer needed from the context, and add information likely to be needed later
- Do this by:
  - Adding an explicit context layer to the architecture
  - This layer controls the flow of information into and out of network layers using specialized neural units called gates

# 

### What does this process look like in a single LSTM unit?












### Gated Recurrent Units (GRUs)

- Also manage the context that is passed through to the next timestep, but do so using a simpler architecture than LSTMs
  - No separate context vector
  - Only two gates
    - Reset gate
    - Update gate
- Gates still use a similar design to that seen in LSTMs
  - Feedforward layer + sigmoid activation + pointwise multiplication with the layer being gated, resulting in a binary-like mask











Review: Essentials of Deep Learning Feedforward neural networks

Convolutional neural networks

Recurrent neural networks

Stacked models

Bidirectional models

### **Stacked RNNs**



- Use the entire sequence of outputs from one RNN as the input sequence to another
- Capable of outperforming single-layer networks
- Why?
  - Having more layers allows the network to learn representations at differing levels of abstraction across layers
    - Early layers → more fundamental properties
    - Later layers → more meaningful groups of fundamental properties
- Optimal number of RNNs to stack together?
  - Depends on application and training set
- More RNNs in the stack  $\rightarrow$  increased training costs

Review: Essentials of Deep Learning Feedforward neural networks

Convolutional neural networks

Recurrent neural networks

Stacked models

**Bidirectional models** 

#### Bidirectional RNNs

- How can we make use of information both before and after the current timestep?
  - Train an RNN on an input sequence in reverse
    - $h_t^b = RNN_{backward}(x_t^n)$ 
      - *h*<sup>b</sup><sub>t</sub> corresponds to information from the current timestep to the end of the sequence
  - Combine the forward and backward networks

### Bidirectional RNNs

- Two independent RNNs
  - One where the input is processed from start to end
  - One where the input is processed from end to start
- Outputs combined into a single representation that captures both the prior and future contexts of an input at each timestep
  - $h_t = h_t^f \oplus h_t^b$
- How to combine the contexts?
  - Concatenation
  - Element-wise addition, multiplication, etc.

### **Bidirectional RNNs**



### **Bidirectional RNNs**



### **Bidirectional RNNs**



#### **Sequence Classification with a Bidirectional RNN**



Natalie Parde - UIC CS 521

## Machine Translation: The process of automatically converting a text from one language to another.



Natalie Parde - UIC CS 521



Le Monde elemondefr Ligue 1 : Lyon rebondit, Angers prend la deuxième place Translated from French by Google

Ligue 1: Lyon bounces back, Angers takes second place



Ligue 1 : Lyon rebondit, Angers prend la deuxième place Face à Nimes, dernier de Ligue 1, les Angevins s'imposent 1-0 et prennent la place de dauphins du PSG. Strasbourg remporte sa première victoire hors de ... & lemonde.fr

81,590 likes maseeloure. If Autumn colors at the taileries guident with thank you (gnathparts for this beautiful perspective) to an average you too, share your most beautiful photos by mentioning #MaseeDuLoure, we share Machine translation is increasingly ubiquitous, but also challenging for many reasons.

Structural and lexical differences between languages

Differences in word order

Morphological differences

Stylistic and cultural differences

#### **Cross-Linguistic Similarities** and Differences

- Typological Differences:
  - Systematic structural differences between languages
- Morphological Differences:
  - Number of morphemes per word
    - Isolating languages: Each word generally has one morpheme
    - Polysynthetic languages: Each word may have many morphemes
  - Degree to which morphemes can be segmented
    - Agglutinative languages: Morphemes have well-defined boundaries
    - Fusion languages: Morphemes may be conflated with one another



### Cross-Linguistic Similarities and Differences

#### Syntactic Differences:

- Primary difference between languages: Word order
  - SVO languages: Verb tends to come between the subject and object
  - SOV languages: Verb tends to come at the end of basic clauses
  - VSO languages: Verb tends to come at the beginning of basic clauses
- Languages with similar basic word order also tend to share other similarities
  - SVO languages generally have prepositions
  - SOV languages generally have postpositions



The bottle exited floating.

### Cross-Linguistic Similarities and Differences

#### Differences in Permissible Omissions:

- Pro-Drop languages: Can omit pronouns when talking about certain referents
- Some pro-drop languages permit more pronoun omission than others
  - Referentially dense and sparse languages
- Converting text from pro-drop languages (e.g., Japanese) to non-pro-drop languages (e.g., English) requires that all missing pronoun locations are identified and their appropriate anaphors recovered
- Differences in noun-adjective order
  - Blue house  $\rightarrow$  Maison bleue
- Differences in homonymy and polysemy
- Differences in grammatical constraints
  - Some languages require gender for nouns
  - Some languages require gender for pronouns
- Lexical gaps
  - No word or phrase in the target language can express the meaning of a word in the source language

### Machine Translation

- Classical Machine Translation
  - Direct translation
  - Transfer approaches
  - Interlingua approaches
  - Statistical methods
- Modern Machine Translation
  - Encoder-decoder models



### Classical Machine Translation

#### Direct translation

- 1. Take a large bilingual dictionary
- 2. Proceed through the source text word by word
- 3. Translate each word according to the dictionary
- No intermediate structures
- Simple reordering rules may be applied
  - For example, moving adjectives so that they are after nouns when translating from English to French
- Dictionary entries may be relatively complex
  - Rule-based programs for translating a word to the target language



## **Classical Machine Translation**

#### Transfer approaches

- Parse the input text
- Apply rules to transform the source language parse structure into a target language parse structure
- Two subcategories of transformations:
  - Syntactic transfer
  - Lexical transfer



#### **Transfer Approaches**





- Syntactic Transfer: Modifies the source parse tree to resemble the target parse tree
  - For some languages, may also include
    thematic structures
    - Directional or locative prepositional phrases vs. recipient prepositional phrases
- Lexical Transfer: Generally based on a bilingual dictionary
  - As with direct translation, dictionary entries can be complex to accommodate many possible translations

#### Classical Machine Translation

#### Interlingua approaches

- Convert the source language text into an abstract meaning representation
- Generate the target language text based on the abstract meaning representation
- Require more analysis work than transfer approaches
  - Semantic analysis
  - Sentiment analysis
- No need for syntactic or lexical transformations



### Interlingua Approaches

- Goal: Represent all sentences that mean the same thing in the same way, regardless of language
- What kind of representation scheme should be used?
  - Classical approaches:
    - First-order logic
    - Semantic primitives
    - Event-based representation
  - More recently, neural machine translation models follow a similar intuition



#### When to use each classical approach?

#### **Direct Translation**

- Pros:
  - Simple
  - · Easy to implement
- Cons:
  - · Cannot reliably handle long-distance reorderings
  - Cannot handle reorderings involving phrases or larger structures
  - · Too focused on individual words

#### Transfer Approaches

- Pros:
  - · Can handle more complex language phenomena than direct translation
- Cons:
  - Still not sufficient for many cases!

#### Interlingua Approaches

- Pros:
  - Direct mapping between meaning representation and lexical realization
- No need for transformation rules
- Cons:
- Extra (often unnecessary) work
- Classical approaches require an exhaustive analysis and formalization of the semantics of the domain

#### Statistical Machine Translation

- Models automatically learn to map from the source language to the target language
  - No need for intermediate transformation rules
  - No need for an explicitly defined internal meaning representation
- Goal: Produce an output that maximizes some function representing translation faithfulness and fluency
- One possible approach: Bayesian noisy channel model
  - Assume a possible target language translation  $t_i$  and a source language sentence s
  - Select the translation t' from the set of all possible translations  $t_i \in T$  that maximizes the probability  $P(t_i|s)$ , using Bayes' rule

• 
$$t' = \underset{t_i \in T}{\operatorname{argmax}} P(s|t_i)P(t_i)$$

• Often a phrase-based translation model is used

#### • • • • • • • • • • • •

## The Phrase-Based **Translation** Model

- Computes the probability that a given translation t<sub>i</sub> generates the original sentence s based on its constituent phrases
- Stages of phrase-based translation:
  - 1. Group the words from the source sentence into phrases
  - 2. Translate each source phrase into a target language phrase
  - 3. (Optionally) reorder the target language phrases

#### **Probability in Phrase-Based Translation Models**

- Relies on two probabilities:
  - Translation probability
    - Probability of generating a source language phrase from a target language phrase,  $\phi(\overline{t_i}, \overline{s_i})$
  - Distortion probability
    - Probability of two consecutive target language phrases being separated in the source language by a word span of a particular length,  $d(a_i b_{i-1})$
- To learn these probabilities, we need to train two sets of parameters:
  - $\phi(\overline{t_i}, \overline{s_i})$
  - $d(a_i b_{i-1})$
- We learn these using phrase-aligned bilingual training sets

Decoding for Phrase-Based Machine Translation

- Aligned phrases can be stored in a phrase-translation table
- **Decoding algorithms** can then search through this table to find the overall translation that maximizes the phrase translation probabilities


Summary: **Review of Deep Learning Architectures** and Classical Machine **Translation** 

- Popular deep learning architectures in natural language processing include feedforward neural networks, convolutional neural networks, and recurrent neural networks
- These architectures can be **stacked** or combined to form
   **bidirectional** architectures
- Machine translation is challenging due to many typological, morphological, and other differences between languages
- Classical machine translation used dictionary-based, direct transfer, and interlingua approaches
- A popular statistical MT model is the Bayesian noisy channel approach, which relies on phrase-based translation

# Machine Translation

- Classical Machine Translation
  - Direct translation
  - Transfer approaches
  - Interlingua approaches
  - Statistical methods
- Modern Machine Translation
  - Encoder-decoder models



#### **Encoder-Decoder Models**

- Generate contextually-appropriate, arbitrary-length output sequences
- Basic premise:
  - Use a neural network to encode an input to an internal representation
  - Pass that internal representation as input to a second neural network
  - Use that neural network to decode the internal representation to a taskspecific output sequence
- This method allows networks to be trained in an end-to-end fashion

# Where did this idea come from?

Sequence processing models generate language using a process of autoregressive generation:

- Start with a seed token (e.g., <s>)
- Predict the most likely next word in the sequence
- Use that word as input at the next timestep
- Repeat until an end token (or max length) is reached



#### This setup can be extended to generate text given a specific prefix....

- Pass the specified prefix through the language model, in sequence
- End with the hidden state corresponding to the last word of the prefix
- Start the autoregressive process at that point
- Goal: Output sequence should be a reasonable completion of the prefix



# We can build upon this idea to transform one type of sequence to another.

• Machine translation example:

- 1. Take a sequence of text from Language #1
- 2. Take the translation of that text from Language #2
- 3. Concatenate the two sequences, separated by a marker
- 4. Use these concatenated sequences to train the autoregressive model
- 5. Test the model by **passing in only the first part of a concatenated sequence** (text from Language #1) and checking to see what the generated completion (text from Language #2) looks like









Key Elements of an Encoder-Decoder Network

#### Encoder

- Accepts an input sequence,  $x_1^n$
- Generates a sequence of contextualized representations,  $h_1^n$

#### Context vector

- A function, c, of h<sub>1</sub><sup>n</sup> that conveys the basic meaning of x<sub>1</sub><sup>n</sup> to the decoder
- (Might just be equivalent to  $h_1^n$ )

#### • Decoder

- Accepts *c* as input
- Generates an arbitrary-length sequence of hidden states,  $h_1^m$ , from which a corresponding sequence of output states  $y_1^m$  can be obtained

# Encoders

- Can be any type of neural network, but is generally assumed to be a sequence processing model:
  - Feedforward network
  - CNN
  - RNN
  - LSTM/BiLSTM
  - GRU/BiGRU
  - Transformer
- More common
- These networks can be stacked on top of one another

# Decoders

- Need to perform autoregressive generation to produce the output sequence
- Can be any type of sequence processing network
  - RNN
  - LSTM
  - GRU
  - Transformer



• 
$$z_t = f(h_t^d)$$

•  $y_t = \operatorname{softmax}(z_t)$ 

# Decoders

- Need to perform autoregressive generation to produce the output sequence
- Can be any type of sequence processing network
  - RNN
  - LSTM
  - GRU
  - Transformer



# Decoders

- Need to perform autoregressive generation to produce the output sequence
- Can be any type of sequence processing network
  - RNN
  - LSTM
  - GRU
  - Transformer
- Formally....
  - $c = h_n^e$
  - $h_0^d = c$

Regular ending steps (activation function applied to hidden state outputs, and softmax applied to activation outputs) •  $h_t^d = g(\widehat{y_{t-1}}, h_{t-1}^d)$ •  $z_t = f(h_t^d)$ •  $y_t = \operatorname{softmax}(z_t)$ 

### A couple useful extensions....

- Formally....
  - $c = h_n^e$
  - $h_0^d = c$
  - $h_t^d = g(\widehat{y_{t-1}}, h_{t-1}^d) \to h_t^d = g(\widehat{y_{t-1}}, h_{t-1}^d, c)$
  - $z_t = f(h_t^d)$
  - $y_t = \operatorname{softmax}(z_t)$

Make the context vector available at each timestep when decoding, so that its influence doesn't diminish over time

### A couple useful extensions....

- Formally....
  - $c = h_n^e$
  - $h_0^d = c$
  - $h_t^d = g(\widehat{y_{t-1}}, h_{t-1}^d) \to h_t^d = g(\widehat{y_{t-1}}, h_{t-1}^d, c)$
  - $z_t = f(h_t^d)$ •  $y_t = \operatorname{softmax}(z_t) \rightarrow y_t = \operatorname{softmax}(\widehat{y_{t-1}}, z_t, c)$

Condition output on not only the hidden state, but the previous output and encoder context (easier to keep track of what's been generated already)

# What other ways can we improve the decoder's output quality?

#### Beam search

- Improved context vector
  - Final hidden state tends to be more focused on the end of the input sequence
  - Can be addressed by using bidirectional RNNs, summing the encoder hidden states, or averaging the encoder hidden states



# **Beam Search**

- Selects from multiple possible outputs by framing the task as a state space search
- Combines breadth-first search with a heuristic filter
  - Continually prunes search space to stay a fixed size (beam width)
- Results in a set of *b* hypotheses, where *b* is the beam width























# How do we choose a best hypothesis?

- Probabilistic scoring scheme
- Pass all or a subset of hypotheses to a downstream application

# So far, the encoder context vectors we've seen have been simple and static.

- Can we do better?
  - Yes!



# Attention Mechanism

- Takes entire encoder context into account
- Can be embodied in a fixed-size vector

### Recall....

- We've already made our context vector available at each timestep when decoding
  h<sup>d</sup><sub>t</sub> = g(ŷ<sub>t-1</sub>, h<sup>d</sup><sub>t-1</sub>, c)
- The first step in creating our attention mechanism is to update our hidden state such that it is conditioned on an updated context vector with each decoding step

• 
$$h_t^d = g(\widehat{y_{t-1}}, h_{t-1}^d, c_t)$$

How do we dynamically create a new context vector at each step?

• Compute a vector of scores that capture the relevance of each encoder hidden state to the decoder hidden state,  $h_{i-1}^d$ 

$$p \ score(h_{i-1}^{a}, h_{j}^{e}) = h_{i-1}^{a} \cdot h_{j}^{e}$$

# **Vector of Context Scores**



# **Vector of Context Scores**



# **Vector of Context Scores**




#### **Vector of Context Scores**



#### **Vector of Context Scores**



#### **Vector of Context Scores**



#### How can we make use of context scores?

- Parameterize these scores with weights
- This allows the model to learn which aspects of similarity between the encoder and decoder states are important

## Attention Weights

- Normalize context scores to create a vector of weights, α<sub>ij</sub>
  - $\alpha_{ij} = \operatorname{softmax}(\operatorname{score}(h_{i-1}^d, h_j^e) \forall j \in e)$
  - Provides the proportional relevance of each encoder hidden state *j* to the current decoder state *i*
- Finally, take a weighted average over all the encoder hidden states to create a fixed-length context vector for the current decoder state

• 
$$c_i = \sum_j \alpha_{ij} h_j^e$$

#### Thus, we finally have an encoderdecoder model with attention!



#### Thus, we finally have an encoderdecoder model with attention!



#### Other Attention Weights

- More sophisticated scoring functions can be used as well
- Common: Parameterize the attention score with its own set of trainable weights
  - score $(\mathbf{h}_{i-1}^d, \mathbf{h}_j^e) = \mathbf{h}_{t-1}^d \mathbf{W}_s \mathbf{h}_j^e$
  - Advantage: Allows the encoder and decoder to use vectors with different dimensionality (dotproduct attention requires the encoder and decoder hidden states to have the same dimensionality)

#### **Practical Details for Building MT Systems**

#### Vocabulary

- MT systems typically use a fixed vocabulary generated using byte pair encoding or other wordpiece algorithms
- Vocabulary is usually shared across the source and target languages

#### Corpora

- Parallel corpora with the same content communicated in multiple languages
- Common sources:
  - Government documents for nations with multiple official languages
  - Subtitles for movies and TV shows
- Often, text from the source and target language(s) is aligned at the sentence level

What if we don't have much training data?

 Parallel corpora are difficult to find, especially for lower-resource language pairs

#### • Backtranslation:

- 1. Train an intermediate target-to-source MT system on a small parallel corpus
- 2. Translate additional monolingual data from the target language to the source language using this intermediate system
- 3. Consider this new, synthetic parallel data as additional training data
- 4. Train a source-to-target MT system on the expanded training dataset

How do we evaluate machine translation models?

- Translation quality tends to be very subjective!
- Two common approaches:
  - Human ratings
  - Automated metrics

Evaluating Machine Translation Using Human Ratings

- Typically evaluated along multiple dimensions
- Tend to check for both fluency and adequacy
- Fluency:
  - Clarity
  - Naturalness
  - Style
- Adequacy:
  - Fidelity
  - Informativeness

## Evaluating Machine Translation Using Human Ratings

- How to get quantitative measures of fluency?
  - Ask humans to rate different aspects of fluency along a scale
  - Measure how long it takes humans to read a segment of text
  - Ask humans to guess the identity of the missing word
    - "After such a late night working on my project, it was hard to wake up this \_\_\_\_!"

### **Evaluating Machine Translation Using Human Ratings**

• How to get quantitative measures of adequacy?

- Ask bilingual raters to rate how much information was preserved in the translation
- Ask monolingual raters to do the same, given access to a gold standard reference translation
- Ask raters to answer multiple-choice questions about content present in a translation

Another set of human evaluation metrics considers postediting.

- Ask a human to post-edit or "fix" a translation
- Compute the number of edits required to correct the output to an acceptable level
  - Can be measured via number of word changes, number of keystrokes, amount of time taken, etc.

## Automated Metrics

- Less accurate than human evaluation, but:
  - Useful for iteratively testing system improvements
  - Can be used as an automatic loss function
- Two main families:
  - Character- or word-overlap
  - Embedding similarity

#### **Popular Lexical Overlap Metrics**

#### • BLEU

- Measure of word overlap
- METEOR
  - Measure of word overlap, considering stemming and synonymy
- Character F-Score (chrF)
  - Measure of character n-gram overlap

#### BLEU

- Weighted average of the number of n-gram overlaps with human translations
- Precision-based metric
  - What percentage of words in the candidate translation also occur in the gold standard translation(s)?
- To compute BLEU:
  - Count how many times each n-gram is used in the candidate translation, c(ngram)
  - Clip that amount so that the highest it can be is c<sub>max</sub>(ngram) defined as the maximum number of times it is used in a reference translation
  - Compute precision for each word in the candidate translation:
    - Divide (a) summed across each candidate, the sum of n-gram frequencies for n-grams appearing in both the candidate and reference, by (b) summed across each candidate, the sum of n-gram frequences for n-grams appearing in the candidate
  - Take the geometric mean of the modified n-gram precisions for unigrams, bigrams, trigrams, and 4-grams

## Then, add a penalty for translation brevity....

- Otherwise, extremely short translations (e.g., "the") could receive perfect scores!
- The penalty is based on:

- The sum of the lengths of the reference sentences, r
- The sum of the lengths of the candidate translations, c
- Formally, the penalty is set to:

• 
$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{(1 - \frac{r}{c})} & \text{if } c \le r \end{cases}$$

- The full BLEU score for a set of translations is then:
  - $BLEU = BP * (\prod_{n=1}^{4} \operatorname{prec}_n)^{\frac{1}{4}}$

Mina no dió una bofetada a la bruja verde.

Source Sentence

Mina didn't slap the green witch.

Reference Translation

Mina did not give a slap to the green witch.

**Candidate Translation** 

Source Sentence

Mina didn't slap the green witch.

Reference Translation

Mina did not give a slap to the green witch.

Candidate Translation

 $prec_n = \frac{\sum_{c \in \{Candidates\}} \sum_{ngram \in c} \min(c(ngram), c_{max}(ngram))}{\sum_{c \in \{Candidates\}} \sum_{ngram \in c} c(ngram)}$ 





Mina didn't slap the green witch.

proc –	$\sum_{c \in \{\text{Candidates}\}} \sum_{\text{ngram} \in C} \min(c(\text{ngram}), c_{\max}(\text{ngram}))$			
$\operatorname{prec}_n$ –	$\sum_{c \in \{Candidates\}} \sum_{ngram \in C} c(ngram)$			
Unigram	Unigram Frequency (Candidate)	Unigram Frequency (Reference)		
Mina	1	1		
did	1	0		
not	1	0		
give	1	0		
а	1	0		
slap	1	1		
to	1	0		
the	1	1		
green	1	1		
witch	1	1		
	1	1		



$$BLEU = BP * (\prod_{n=1}^{4} \operatorname{prec}_{n})^{\frac{1}{4}}$$

Mina didn't slap the green witch.

prog -	$\sum_{c \in \{Candidates\}} \sum_{ngram \in c} \min(c(ngram), c_{max}(ngram))$			
$\Sigma_{c \in \{\text{Candidates}\}} \Sigma_{n \text{gram} \in C} c(n \text{gram})$				
Unigram	Unigram Frequency (Candidate)	Unigram Frequency (Reference)		
Mina	1	1		
did	1	0		
not	1	0		
give	1	0	1 +	
а	1	0	$p_1 = \frac{1}{1}$	
slap	1	1	1 +	
to	1	0		
the	1	1		
green	1	1		
witch	1	1		
	1	1		



$$BLEU = BP * (\prod_{n=1}^{4} \operatorname{prec}_{n})^{\frac{1}{4}}$$

$$p_1 = \frac{1+0+0+0+0+1+0+1+1+1+1}{1+1+1+1+1+1+1+1+1} = \frac{6}{11}$$

Mina didn't slap the green witch.

$$prec_n = \frac{\sum_{c \in \{Candidates\}} \sum_{ngram \in C} \min(c(ngram), c_{max}(ngram))}{\sum_{c \in \{Candidates\}} \sum_{ngram \in C} c(ngram)}$$



$$p_1 = \frac{1+0+0+0+0+1+0+1+1+1+1}{1+1+1+1+1+1+1+1+1} = \frac{6}{11}$$

$$p_2 = \frac{0+0+0+0+0+0+0+1+1+1}{1+1+1+1+1+1+1+1} = \frac{3}{10}$$

Bigram	Bigram Frequency (Candidate)	Bigram Frequency (Reference)
Mina did	1	0
did not	1	0
not give	1	0
give a	1	0
a slap	1	0
slap to	1	0
to the	1	0
the green	1	1
green witch	1	1
witch.	1	1

**Trigram Frequency** 

Mina didn't slap the green witch.

Trigram

Mina did not give a slap to the green witch.

$$\operatorname{prec}_{n} = \frac{\sum_{c \in \{\text{Candidates}\}} \sum_{\operatorname{ngram} \in C} \min(c(\operatorname{ngram}), c_{\max}(\operatorname{ngram}))}{\sum_{c \in \{\text{Candidates}\}} \sum_{\operatorname{ngram} \in C} c(\operatorname{ngram})}$$

$$BP = \begin{cases} 1 & if \ c > r \\ e^{(1 - \frac{r}{c})} & if \ c \le r \end{cases}$$

$$BLEU = BP * (\prod_{n=1}^{4} \operatorname{prec}_{n})^{\frac{1}{4}}$$

**Trigram Frequency** 

$$p_1 = \frac{6}{11} \qquad p_2 = \frac{3}{10}$$

$$p_3 = \frac{0+0+0+0+0+0+0+1+1}{1+1+1+1+1+1+1+1} = \frac{2}{9}$$

Mina didn't slap the green witch.

$$\operatorname{prec}_{n} = \frac{\sum_{c \in \{\text{Candidates}\}} \sum_{\operatorname{ngram} \in C} \min(c(\operatorname{ngram}), c_{\max}(\operatorname{ngram}))}{\sum_{c \in \{\text{Candidates}\}} \sum_{\operatorname{ngram} \in C} c(\operatorname{ngram})}$$

$$BP = \begin{cases} 1 & if \ c > r \\ e^{(1 - \frac{r}{c})} \ if \ c \le r \end{cases}$$

$$BLEU = BP * (\prod_{n=1}^{4} \operatorname{prec}_{n})^{\frac{1}{4}}$$

$$p_1 = \frac{6}{11} \qquad p_2 = \frac{3}{10} \qquad p_3 = \frac{2}{9}$$
$$p_4 = \frac{0+0+0+0+0+0+0+1}{1+1+1+1+1+1+1} = \frac{1}{8}$$

4-gram	4-gram Frequency (Candidate)	4-gram Frequency (Reference)
Mina did not give	1	0
did not give a	1	0
not give a slap	1	0
give a slap to	1	0
a slap to the	1	0
slap to the green	1	0
to the green witch	1	0
the green witch .	1	1

Mina didn't slap the green witch.

Mina did not give a slap to the green witch.

 $\operatorname{prec}_{n} = \frac{\sum_{c \in \{\operatorname{Candidates}\}} \sum_{\operatorname{ngram} \in C} \min(\operatorname{c(ngram)}, \operatorname{c_{max}(ngram)})}{\sum_{c \in \{\operatorname{Candidates}\}} \sum_{\operatorname{ngram} \in C} \operatorname{c(ngram)}}$ 



*c* = 11

$$BLEU = BP * (\prod_{n=1}^{4} \operatorname{prec}_{n})^{\frac{1}{4}}$$

$$p_1 = \frac{6}{11}$$
  $p_2 = \frac{3}{10}$   $p_3 = \frac{2}{9}$   $p_4 = \frac{1}{8}$   
 $BP = 1$ 

 $BLEU = 1 * \left(\prod_{n=1}^{4} \operatorname{prec}_{n}\right)^{\frac{1}{4}} = 1 * \left(\frac{6}{11} * \frac{3}{10} * \frac{2}{9} * \frac{1}{8}\right)^{\frac{1}{4}} = 1 * 0.0045454545454^{\frac{1}{4}} = 1 * 0.25965358893 = 0.26$ 

# What are<br/>good<br/>BLEUScores?

+

0

- No formal score ranges, but in general:
  - < 10: Very poor
  - 10-19: Difficult to understand
  - 20-29: Understandable but with many grammatical errors
  - 30-39: Understandable and reasonable quality
  - 40-49: High quality
  - 50-59: Very high quality
  - >60: May exceed human translators

#### **Limitations of BLEU**

- Word or phrase order is of minimal importance
  - When computing unigram precision, a word can exist anywhere in the translation!
- Does not consider word similarity
- Relatively low correlation with human ratings
- Nonetheless, BLEU is reasonable to use in cases when a quick, automated metric is needed to assess translation performance

## Character F-Score (chrF)

+

0

- Same intuition as BLEU: Good machine translations tend to contain the same words and characters as human translations
- Ranks a candidate translation by a function of the number of character ngram overlaps with a human reference translation
- Less sensitive to differences in wordform than BLEU (e.g., "gives" versus "is giving")

## How is chrF computed?

- Similarly to "regular" F-score
  - chrP: averaged % of character unigrams, bigrams, ..., k-grams in the hypothesis that o also occur in the reference
  - **chrR:** averaged % of character unigrams, bigrams, ..., k-grams in the reference that also occur in the hypothesis
  - β: weighting parameter (similarly to F-score) that determines the relative impacts of chrP and chrR on the overall F-score

• 
$$\operatorname{chr}F\beta = (1 + \beta^2) \frac{\operatorname{chr}P \times \operatorname{chr}R}{\beta^2 \times \operatorname{chr}P + \operatorname{chr}R}$$
, or for example,  $\operatorname{chr}F2 = \frac{5 \times \operatorname{chr}P \times \operatorname{chr}R}{4 \times \operatorname{chr}P + \operatorname{chr}R}$ 

+















**chrR:** averaged % of character unigrams, bigrams, ..., k-grams in the reference that are also in the hypothesis

Natalie Parde - UIC CS 521




9 unigrams are in both the hypothesis and the reference





9 unigrams are in both the hypothesis and the reference

6 bigrams are in both the hypothesis and the reference





9 unigrams are in both the hypothesis and the reference

6 bigrams are in both the hypothesis and the reference

5 trigrams are in both the hypothesis and the reference





9 unigrams are in both the hypothesis and the reference

Unigram chrP:  $^{9}/_{12} = 0.75$ 

6 bigrams are in both the hypothesis and the reference

5 trigrams are in both the hypothesis and the reference





5 trigrams are in both the hypothesis and the reference









Bigram chrR:  $^{6}/_{13} = 0.46$ 

Trigram chrR:  $\frac{5}{12} = 0.42$ 

Bigram chrP: 
$$^{6}/_{11} = 0.55$$

Trigram chrP: 
$$\frac{5}{10} = 0.5$$

chrP: 
$$\frac{0.75+0.55+0.5}{3} = 0.6$$

Natalie Parde - UIC CS 521





Bigram chrR:  $\frac{6}{13} = 0.46$ 

Trigram chrR:  $\frac{5}{12} = 0.42$ 

chrR: 0.64+0.46+0.42

Bigram chrP: 
$$^{6}/_{11} = 0.55$$

Trigram chrP: 
$$\frac{5}{10} = 0.5$$
  
chrP:  $\frac{0.75 + 0.55 + 0.5}{10} = 0.6$ 

Natalie Parde - UIC CS 521

 $\frac{2}{2} = 0.51$ 





chrP: 
$$\frac{0.75+0.55+0.5}{3} = 0.6$$

chrR: 
$$\frac{0.64+0.46+0.42}{3} = 0.51$$

$$chrF2 = \frac{5 * chrP * chrR}{4 * chrP + chrR} = \frac{5 * 0.6 * 0.51}{4 * 0.6 + 0.51} = 0.53$$

# Limitations of chrF

- Focuses on differences at a very local scale (i.e., character n-grams)
- Doesn't measure discourse coherence
- Best at measuring performance for different versions of the same system, rather than comparing different systems

## Embedding -Based Evaluation Methods

- Measuring exact word- or character-level overlap might be overly strict
  - Good translations may use words that are synonymous to those in the reference!
- Embedding-based methods measure the semantic overlap between reference and hypothesis translations

#### Popular Embedding-Based Methods for Evaluating MT Systems



### **Summary: Encoder**-Decoder **Models** and Evaluating **MT Systems**

- Encoder-decoder models draw upon similar techniques for autoregressive language modeling to convert input to an intermediate vector representation and then convert that intermediate representation to output
- Attention scores help the encoderdecoder model's context vector focus on the most relevant information from the input for a given decoder timestep
- MT systems are commonly evaluated using both human ratings and automated metrics
- Popular automated metrics include BLEU, chrF, and embedding-based measures