Human-Inspired Structured Prediction for Language and Biology

Assistant Professor, Oregon State University

Liang Huang



Principal Scientist, Baidu Research

incremental & linear-time Human-Inspired Structured Prediction for Language and Biology

Assistant Professor, Oregon State University

Liang Huang



Principal Scientist, Baidu Research

incremental & linear-time Human-Inspired Structured Prediction for Language and Biology

simultaneous interpretation





Liang Huang

Assistant Professor, Oregon State University



Principal Scientist, Baidu Research

incremental & linear-time Human-Inspired Structured Prediction for Language and Biology

simultaneous interpretation



syntactic structure





Principal Scientist, Baidu Research Assistant Professor, Oregon State University

Liang Huang





incremental & linear-time

simultaneous interpretation





natural language sequence



RNA sequence

Liang Huang



- Principal Scientist, Baidu Research
- Assistant Professor, Oregon State University

My PhD Graduates





Ashish Vaswani (USC, 2014) (co-advised by David Chiang)

Senior Research Scientist Google Brain

first author of Transformer "Attention is All You Need" James Cross (OSU, 2016)

Research Scientist Facebook

EMNLP 2016 Best Paper Honorable Mention

A Bit about Myself...





Kai Zhao (OSU, 2017)

Research Scientist Google

11 top-conference papers (ACL/EMNLP/NAACL)

Mingbo Ma (OSU, 2018)

Research Scientist Baidu Research USA

breakthrough in simultaneous translation



My PhD Graduates



Ashish Vaswani (USC, 2014) (co-advised by David Chiang)

Senior Research Scientist Google Brain

first author of Transformer "Attention is All You Need"



James Cross (OSU, 2016) **Kai Zhao** (OSU, 2017)

Research Scientist Facebook

EMNLP 2016 Best Paper Honorable Mention

My Research: Efficient Structured Prediction

Bush met Putin in Moscow

source language sentence

l eat sushi with tuna from Japan GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA

natural language sentence

A Bit about Myself...





Research Scientist Google

11 top-conference papers (ACL/EMNLP/NAACL)

Mingbo Ma (OSU, 2018)

Research Scientist Baidu Research USA

breakthrough in simultaneous translation

RNA sequence







My PhD Graduates



Ashish Vaswani (USC, 2014) (co-advised by David Chiang)

Senior Research Scientist Google Brain

first author of Transformer "Attention is All You Need"

target-language sequence

布什在莫斯科与普京会晤

Bush met Putin in Moscow

source language sentence



James Cross (OSU, 2016)

Research Scientist Facebook

EMNLP 2016 Best Paper Honorable Mention

syntactic structure



natural language sentence

A Bit about Myself...





Kai Zhao (OSU, 2017)

Research Scientist Google

11 top-conference papers (ACL/EMNLP/NAACL)

Mingbo Ma (OSU, 2018)

Research Scientist Baidu Research USA

breakthrough in simultaneous translation

My Research: Efficient Structured Prediction

secondary structure



• how many interpretations?

I saw her duck







• how many interpretations?

I saw her duck













• how many interpretations?

I eat sushi with tuna





• how many interpretations?



I eat sushi with tuna



www.shutterstock.com · 46863619





• how many interpretations?







structural ambiguity



• how many interpretations?







structural ambiguity







Unexpected Structural Ambiguity







• how many interpretations?

I saw her duck

















• how many interpretations?

But humans can resolve these ambiguities incremental in linear-time!

I saw her duck with a telescope in the garden ...











I eat sushi with tuna from Japan





• human sentence processing is well-known to be incremental and linear-time





I eat sushi with tuna from Japan





• human sentence processing is well-known to be incremental and linear-time











• human sentence processing is well-known to be incremental and linear-time





they often need full sentence as input and run in superlinear time





natural language processing algorithms are often non-incremental and slow



• human sentence processing is well-known to be incremental and linear-time



they often need full sentence as input and run in superlinear time









natural language processing algorithms are often non-incremental and slow



• human sentence processing is well-known to be incremental and linear-time



they often need full sentence as input and run in superlinear time









natural language processing algorithms are often non-incremental and slow





Three Stories on Incrementality and Instantaneity

simultaneous translation incremental parsing **linear-time RNA structure prediction**





Three Stories on Incrementality and Instantaneity

linear-time RNA structure prediction simultaneous translation incremental parsing

Bush met Putin in Moscow

source language sentence

l eat sushi with tuna from Japan GCGGGAAUAGCUCAGUUGGUAGAGCACGACCU

natural language sentence

RNA sequence





8

Three Stories on Incrementality and Instantaneity

linear-time RNA structure prediction simultaneous translation incremental parsing

target-language sequence

布什在莫斯科与普京会晤

Bush met Putin in Moscow

source language sentence

syntactic structure



natural language sentence

secondary structure





Part I: Simultaneous Translation

(Ma, Huang, et al, ArXiv 2018; under review)

Part I: Simultaneous Translation

(Ma, Huang, et al, ArXiv 2018; under review)

Background: Consecutive vs. Simultaneous Interpretation

consecutive interpretation *multiplicative latency (x2)*



simultaneous interpretation additive latency (+3 secs)







Background: Consecutive vs. Simultaneous Interpretation

consecutive interpretation *multiplicative latency (x2)*



simultaneous interpretation additive latency (+3 secs)

simultaneous interpretation is extremely difficult

only ~3,000 qualified simultaneous interpreters world-wide



each interpreter can only sustain for at most 10-30 minutes

the best interpreters can only cover ~60% of the source material







Background: Consecutive vs. Simultaneous Interpretation

consecutive interpretation *multiplicative latency (x2)*



just use standard full-sentence translation (e.g., seq-to-seq)

simultaneous interpretation additive latency (+3 secs)

simultaneous interpretation is extremely difficult

only ~3,000 qualified simultaneous interpreters world-wide



each interpreter can only sustain for at most 10-30 minutes

the best interpreters can only cover ~60% of the source material

one of the holy grails of AI

need fundamentally different ideas!





Our Breakthrough

our

full-sentence (non-simultaneous) translation latency: one sentence (10+ secs)



Baidu World Conference, November 2017

and many other companies

simultaneous translation achieved for the first time latency ~3 secs



Baidu World Conference, November 2018


our

full-sentence (non-simultaneous) translation latency: one sentence (10+ secs)



Baidu World Conference, November 2017

and many other companies

simultaneous translation achieved for the first time latency ~3 secs



Baidu World Conference, November 2018



our

full-sentence (non-simultaneous) translation latency: one sentence (10+ secs)



Baidu World Conference, November 2017

and many other companies

simultaneous translation achieved for the first time latency ~3 secs



Baidu World Conference, November 2018



our

full-sentence (non-simultaneous) translation latency: one sentence (10+ secs)



Baidu World Conference, November 2017

and many other companies

simultaneous translation achieved for the first time latency ~3 secs



Baidu World Conference, November 2018



our

full-sentence (non-simultaneous) translation latency: one sentence (10+ secs)



Baidu World Conference, November 2017

and many other companies

simultaneous translation achieved for the first time latency ~3 secs



Baidu World Conference, November 2018



- e.g. translate from Subj-Obj-Verb (Japanese, German) to Subj-Verb-Obj (English)
 - German is underlyingly SOV, and Chinese is a mix of SVO and SOV
 - human simultaneous interpreters routinely "anticipate" (e.g., predicting German verb)

ich bin mit dem Zug nach Ulm gefahren am with the train to Ulm **traveled**

Grissom et al, 2014

 $(\ldots waiting.\ldots)$ **traveled** by train to Ulm



- e.g. translate from Subj-Obj-Verb (Japanese, German) to Subj-Verb-Obj (English)
 - German is underlyingly SOV, and Chinese is a mix of SVO and SOV
 - human simultaneous interpreters routinely "anticipate" (e.g., predicting German verb)

ich bin mit dem Zug nach Ulm gefahren am with the train to Ulm **traveled** zŏngtŏng Bùshí Mòsīkē zài уŭ 总统 在 与 俄罗斯 莫斯科 布什 President with Russian Bush Moscow in

President Bush meets with Russian President Putin in Moscow

Grissom et al, 2014

 $(\ldots waiting.\ldots)$ **traveled** by train to Ulm





- e.g. translate from Subj-Obj-Verb (Japanese, German) to Subj-Verb-Obj (English)
 - German is underlyingly SOV, and Chinese is a mix of SVO and SOV
 - human simultaneous interpreters routinely "anticipate" (e.g., predicting German verb)
- ich bin mit dem Zug nach Ulm gefahren am with the train to Ulm **traveled** $(\ldots waiting.\ldots)$ **traveled** by train to Ulm zŏngtŏng Éluósī Bùshí Mòsīkē zŏngtŏng Pŭjīng zài huìwù уŭ 总统 在 与 俄罗斯 总统 莫斯科 布什

with

President Bush meets with Russian President Putin in Moscow non-anticipative: President Bush (..... waiting)

Moscow

President

in

Bush

Grissom et al, 2014









meets with Russian ...





- e.g. translate from Subj-Obj-Verb (Japanese, German) to Subj-Verb-Obj (English)
 - German is underlyingly SOV, and Chinese is a mix of SVO and SOV
 - human simultaneous interpreters routinely "anticipate" (e.g., predicting German verb)
- ich bin mit dem Zug nach Ulm gefahren am with the train to Ulm **traveled** $(\ldots waiting.\ldots)$ **traveled** by train to Ulm



President Bush meets with Russian President Putin in Moscow non-anticipative: President Bush (..... waiting)

Grissom et al, 2014







- meets with Russian ...
- anticipative: President Bush meets with Russian President Putin in Moscow





- standard seq-to-seq is only suitable for conventional full-sentence MT
- we propose prefix-to-prefix, tailed to simultaneous MT
 - special case: wait-k policy: translation is always k words behind source sentence
 - training in this way enables anticipation



- standard seq-to-seq is only suitable for conventional full-sentence MT
- we propose prefix-to-prefix, tailed to simultaneous MT
 - special case: wait-k policy: translation is always k words behind source sentence
 - training in this way enables anticipation







- standard seq-to-seq is only suitable for conventional full-sentence MT
- we propose prefix-to-prefix, tailed to simultaneous MT
 - special case: wait-k policy: translation is always k words behind source sentence
 - training in this way enables anticipation







- standard seq-to-seq is only suitable for conventional full-sentence MT
- we propose prefix-to-prefix, tailed to simultaneous MT
 - special case: wait-k policy: translation is always k words behind source sentence
 - training in this way enables anticipation

President Bush

meets



wait 2



- standard seq-to-seq is only suitable for conventional full-sentence MT
- we propose prefix-to-prefix, tailed to simultaneous MT
 - special case: wait-k policy: translation is always k words behind source sentence
 - training in this way enables anticipation



wait 2

President Bush meets with



- standard seq-to-seq is only suitable for conventional full-sentence MT
- we propose prefix-to-prefix, tailed to simultaneous MT
 - special case: wait-k policy: translation is always k words behind source sentence
 - training in this way enables anticipation



wait 2



President Bush meets with Russian

- standard seq-to-seq is only suitable for conventional full-sentence MT
- we propose prefix-to-prefix, tailed to simultaneous MT

wait 2

- special case: wait-k policy: translation is always k words behind source sentence
- training in this way enables anticipation





President Bush meets with Russian President

- standard seq-to-seq is only suitable for conventional full-sentence MT
- we propose prefix-to-prefix, tailed to simultaneous MT
 - special case: wait-k policy: translation is always k words behind source sentence
 - training in this way enables anticipation



wait 2



zŏngtŏng Půjīng 总统 Putin President

普京

President Bush meets with Russian President Putin

- standard seq-to-seq is only suitable for conventional full-sentence MT
- we propose prefix-to-prefix, tailed to simultaneous MT

wait 2

- special case: wait-k policy: translation is always k words behind source sentence
- training in this way enables anticipation









President Bush meets with Russian President Putin in Moscow

Research Demo

江泽民对法国总统的来华 jiang zemin expressed his appreciation







Research Demo

江泽民对法国总统的来华 jiang zemin expressed his appreciation







Research Demo

江泽民对法国总统的来华 iang zemin expressed his appreciation

zémín duì fǎ guó zǒng tǒng de láihuá jiāng 江泽民对法国总统 的 来华 访问 to-China visit jiang zemin to French President 's

jiang zemin expressed his appreciation for





Bai Research This is just our research demo. Our production system is better (shorter ASR latency).



Chinese input:	
Pinyin:	
Word-by-Word Translation:	
Simultaneous Translation (wait 3):	
Simultaneous Translation (wait 5):	
Baseline Tranlation (gready):	
Baseline Tranlation (beam 5):	

Latency-Accuracy Tradeoff

Anna Anna Anna Anna Anna	na analas analas sanana analas sanana s	ananan sananan sananan sananan sananan sanan	na vanana vanana vanana vanana vanana vanana	analasi analasi analasi analasi analasi a	anna ann an





Chinese input:	
Pinyin:	
Word-by-Word Translation:	
Simultaneous Translation (wait 3):	
Simultaneous Translation (wait 5):	
Baseline Tranlation (gready):	
Baseline Tranlation (beam 5):	

Latency-Accuracy Tradeoff

Anna Anna Anna Anna Anna	na analas analas sanana analas sanana s	ananan sananan sananan sananan sananan sanan	na vanana vanana vanana vanana vanana vanana	analasi analasi analasi analasi analasi a	anna ann an





Deployment Demo



Bai Research This is live recording from the Baidu World Conference on Nov 1, 2018.





Deployment Demo



Bai Research This is live recording from the Baidu World Conference on Nov 1, 2018.





Experimental Results (German=>English)

German source:

doch während man sich im kongress nicht auf ein vorgehen einigen kann, warten mehrere bundesstaaten nicht länger. they self in congress not on one action agree can wait several states while but not

English translation (simultaneous, wait 3): but, while congress <u>does</u> not agree on a course of action, several states no longer wait.

English translation (full-sentence baseline):

but, while congressional action can not be agreed, several states are no longer waiting.





Experimental Results (German=>English)

German source:

doch während man sich im kongress nicht auf ein vorgehen einigen kann, warten mehrere bundesstaaten nicht länger. they self in congress not on one action agree while wait several but states can not

English translation (simultaneous, wait 3):

English translation (full-sentence baseline): but, while congressional action can not be agreed, several states are no longer waiting.



but, while congress <u>does</u> not agree on a course of action, several states no longer wait.





Experimental Results (German=>English)

German source:

doch während man sich im kongress nicht auf ein vorgehen einigen kann, warten mehrere bundesstaaten nicht länger. they self in congress not on one action agree while several states wait but can not

English translation (simultaneous, wait 3):

English translation (full-sentence baseline): but, while congressional action can not be agreed, several states are no longer waiting.



but, while congress <u>does</u> not agree on a course of action, several states no longer wait.

ich bin mit dem Zug nach Ulm **gefahren** am with the train to Ulm traveled wait 8 words I traveled to Ulm by train full-sentence baseline: CW = 8wait 6 words wait 2 traveled to Ulm by train Gu et al. (2017): CW = (2+6)/2 = 4to Ulm I 1 took train wait 4 a our wait 4 model: CW = (4+1+1+1+1)/5 = 1.6







Summary of Innovations and Impact

- first simultaneous translation approach with integrated anticipation
 - inspired by human simultaneous interpreters who routinely anticipate
- first simultaneous translation approach with arbitrary controllable latency
 - previous RL-based approaches can encourage but can't enforce latency limit
- very easy to train and scalable minor changes to any neural MT codebase
- prefix-to-prefix is very general; can be used in other tasks with simultaneity





Summary of Innovations and Impact

- first simultaneous translation approach with integrated anticipation
 - inspired by human simultaneous interpreters who routinely anticipate
- first simultaneous translation approach with arbitrary controllable latency
 - previous RL-based approaches can encourage but can't enforce latency limit
- very easy to train and scalable minor changes to any neural MT codebase
- prefix-to-prefix is very general; can be used in other tasks with simultaneity



Next: Integrate Incremental Predictive Parsing

• how to be smarter about when to wait and when to translate?

mandatory reordering (i.e., wait):

x í jìnpíng y ú nián zài běijīng dāngxuǎn 习近平于 2012 年 在 北京 当选 Xi Jiping in 2012 yr in Beijing elected

reference translation

"Xi Jinping was elected in Beijing in 2012"

optional reordering:

méiyǒu zhǔnquè guānyú kèlíndùnzhǔyì d e dìngyì 关于克林顿主义, 没有 准确 的 定义 about Clintonism def. accurate no

"There is no accurate definition of Clintonism."







Next: Integrate Incremental Predictive Parsing

• how to be smarter about when to wait and when to translate?

mandatory reordering (i.e., wait):

x í jìnpíng y ú nián zài běijīng dāngxuǎn 习近平于 2012 年 在 北京 当选 Xi Jiping in 2012 yr in Beijing elected

reference translation

"Xi Jinping was elected in Beijing in 2012"

ideal simultaneous

Xi Jinping

was elected...

optional reordering:

méiyǒu zhǔnquè guānyú kèlíndùnzhǔyì d e dìngyì 关于克林顿主义, 没有 准确 的 定义 about Clintonism def. accurate no

"There is no accurate definition of Clintonism."

About Clintonism, there is no accurate definition.









Next: Integrate Incremental Predictive Parsing

• how to be smarter about when to wait and when to translate?





"There is no accurate definition of Clintonism."

About Clintonism, there is no accurate definition.









constituency parsing

Part II: Linear-Time Incremental Parsing

(Huang & Sagae, ACL 2010^{*}; Goldberg, Zhao, Huang, ACL 2013; Zhao, Cross, Huang, EMNLP 2013; Mi & Huang, ACL 2015; Cross & Huang, ACL 2016; Cross & Huang, EMNLP 2016** Hong and Huang, ACL 2018)



dependency parsing

* best paper nominee ** best paper honorable mention



constituency parsing

Part II: Linear-Time Incremental Parsing

(Huang & Sagae, ACL 2010^{*}; Goldberg, Zhao, Huang, ACL 2013; Zhao, Cross, Huang, EMNLP 2013; Mi & Huang, ACL 2015; Cross & Huang, ACL 2016; Cross & Huang, EMNLP 2016** Hong and Huang, ACL 2018)



dependency parsing

* best paper nominee ** best paper honorable mention

Motivations for Incremental Parsing

- simultaneous translation
- auto completion (search suggestions)
- question answering
- dialog
- speech recognition
- input method editor

siri is so

siri is so dumb siri is so bad siri is so **useless** siri is so slow siri is so **rude**







Human Parsing vs. Compilers vs. NL Parsing




Human Parsing vs. Compilers vs. NL Parsing





Human Parsing vs. Compilers vs. NL Parsing



- can we design NL parsing algorithms that is both fast and accurate, inspired by human sentence processing and compilers?
- our idea: generalize PL parsing (LR algorithm) to NL parsing, but keep it O(n)
- challenge: how to deal with ambiguity explosion in NL?
- solution: linear-time dynamic programming both fast and accurate!



Solution: linear-time, DP, and accurate!

- very fast linear-time dynamic programming parser
- explores exponentially many trees (and outputs forest)
- accurate parsing accuracy on English & Chinese







Solution: linear-time, DP, and accurate!

- very fast linear-time dynamic programming parser
- explores exponentially many trees (and outputs forest)
- accurate parsing accuracy on English & Chinese



 $O(n^2)$





Solution: linear-time, DP, and accurate!

- very fast linear-time dynamic programming parser
- explores exponentially many trees (and outputs forest)
- accurate parsing accuracy on English & Chinese



 $O(n^{2.4})$ 10¹⁰ 10⁸ $O(n^2)$ 10⁶ 10^{4} 10² non-DP beam search 10⁰ 10 20 30 40 50 60 70 0 sentence length





I eat sushi with tuna from Japan





I eat sushi with tuna from Japan

Incremental Parsing (Shift-Reduce)





leat sushi with tuna from Japan

stack action

Incremental Parsing (Shift-Reduce)

queue





leat sushi with tuna from Japan

stack action 0

Incremental Parsing (Shift-Reduce)

queue

eat sushi ...







queue

eat sushi ...

eat sushi with ...





queue

eat sushi ...

eat sushi with ...

sushi with tuna ...





queue

eat sushi ...

eat sushi with ...

sushi with tuna ...

sushi with tuna ...





queue

eat sushi ...

eat sushi with ...

sushi with tuna ...

sushi with tuna ...

with tuna from ...





queue

eat sushi ...

eat sushi with ...

sushi with tuna ...

sushi with tuna ...

with tuna from ...

with tuna from ...





0

2

3

4

5a

queue

eat sushi ...

eat sushi with ...

sushi with tuna ...

sushi with tuna ...

with tuna from ...

with tuna from ...

tuna from Japan ...





queue

eat sushi ...

eat sushi with ...

sushi with tuna ...

sushi with tuna ...

with tuna from ...

with tuna from ...

tuna from Japan ...

shift-reduce conflict



Greedy Search

- each state => three new states (shift, l-reduce, r-reduce)
- greedy search: always pick the best next state
 - "best" is defined by a score learned from data







Greedy Search

- each state => three new states (shift, I-reduce, r-reduce)
- greedy search: always pick the best next state

• "best" is defined by a score learned from data







Beam Search

- each state => three new states (shift, I-reduce, r-reduce)
- beam search: always keep top-b states
 - still just a tiny fraction of the whole search space







Beam Search

- each state => three new states (shift, l-reduce, r-reduce)
- beam search: always keep top-b states
 - still just a tiny fraction of the whole search space

psycholinguistic evidence: parallelism (Fodor et al, 1974; Gibson, 1991)



- each state => three new states (shift, l-reduce, r-reduce)
- key idea of DP: share common subproblems

merge equivalent states => polynomial space





- each state => three new states (shift, l-reduce, r-reduce)
- key idea of DP: share common subproblems

merge equivalent states => polynomial space







- each state => three new states (shift, l-reduce, r-reduce)
- key idea of DP: share common subproblems
 - merge equivalent states => polynomial space







- each state => three new states (shift, l-reduce, r-reduce)
- key idea of DP: share common subproblems
 - merge equivalent states => polynomial space

each DP state corresponds to exponentially many non-DP states



graph-structured stack (Tomita, 1986)







- each state => three new states (shift, I-reduce, r-reduce)
- key idea of DP: share common subproblems
 - merge equivalent states => polynomial space

each DP state corresponds to exponentially many non-DP states







- each state => three new states (shift, I-reduce, r-reduce)
- key idea of DP: share common subproblems
 - merge equivalent states => polynomial space

each DP state corresponds to exponentially many non-DP states



graph-structured stack (Tomita, 1986)





- two states are equivalent if they agree on features
 - because same features guarantee same cost
 - example: if we only care about the last 2 words on stack





- two states are equivalent if they agree on features
 - because same features guarantee same cost
 - example: if we only care about the last 2 words on stack





- two states are equivalent if they agree on features
 - because same features guarantee same cost
 - example: if we only care about the last 2 words on stack



psycholinguistic evidence (eye-tracking experiments):

delayed disambiguation

John and Mary had 2 papers John and Mary had 2 papers

Frazier and Rayner (1990), Frazier (1999)



- two states are equivalent if they agree on features
 - because same features guarantee same cost
 - example: if we only care about the last 2 words on stack



psycholinguistic evidence (eye-tracking experiments):

delayed disambiguation

John and Mary had 2 papers each John and Mary had 2 papers together

Frazier and Rayner (1990), Frazier (1999)



Results: Fast and Accurate

Constituency parsing, PTB only, Single Model, End-to-End

Parsre	Note	F1 Score
Durett + Klein 2015	cubic-time parser	91.1
Cross + Huang 2016	original span parser (greedy)	91.3
Liu + Zhang 2016	greedy / beam	91.7
Dyer et al. 2016	greedy / beam	91.7
Stern 2017a	cubic-time span-based parser	91.79
Our Work	linear-time dynamic programming, span-based	91.97



(Hong and Huang, 2018)



Part III: Linear-Time RNA Structure Prediction

(Huang et al, 2019; under review)

Part III: Linear-Time RNA Structure Prediction

(Huang et al, 2019; under review)



Computational Linguistics => Computational Biology





Computational Linguistics => Computational Biology





RNAs and Structure Prediction



GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA






RNAs and Structure Prediction









allowed pairs: G-C A-U G-U assume no crossing pairs

input GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA

example: transfer RNA (tRNA)





allowed pairs: G-C A-U G-U assume no crossing pairs

input \mathcal{X} GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA output

example: transfer RNA (tRNA)





allowed pairs: G-C A-U G-U assume no crossing pairs

> input ${\mathcal X}$ output \mathcal{V}



example: transfer RNA (tRNA)





allowed pairs: G-C A-U G-U assume no crossing pairs

> input $\boldsymbol{\chi}$ output \mathcal{V}









parse tree





allowed pairs: G-C A-U G-U assume no crossing pairs

> input output









allowed pairs: G-C A-U G-U assume no crossing pairs









GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA

- idea 0: tag each nucleotide from left to right maintain a stack: push "(", pop ")", skip "."
 - exhaustive: $O(3^n)$



GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA

- idea 0: tag each nucleotide from left to right
 - maintain a stack: push "(", pop ")", skip "."
 - exhaustive: $O(3^n)$



GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA

- idea I: DP by merging "equivalent states"
 - maintain graph-structured stacks

• DP: $O(n^3)$







GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA

- idea I: DP by merging "equivalent states"
 - maintain graph-structured stacks

• DP: $O(n^3)$







GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA

- idea 2: approximate search: beam pruning
 - keep only top b states per step
 - DP+beam: O(n)







GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA

- idea 2: approximate search: beam pruning
 - keep only top b states per step
 - DP+beam: O(n)

each DP state corresponds to exponentially many non-DP states





GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA

- idea 2: approximate search: beam pruning
 - keep only top b states per step
 - DP+beam: O(n)

each DP state corresponds to exponentially many non-DP states





GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA

- idea 2: approximate search: beam pruning
 - keep only top b states per step
 - DP+beam: O(n)

each DP state corresponds to exponentially many non-DP states



Sequence Length

Ambiguity Packing in Biology and Language

• two states are "temporarily equivalent" if their rightmost unpaired brackets are the same



psycholinguistic evidence (eye-tracking experiments):

delayed disambiguation

John and Mary had 2 papers John and Mary had 2 papers

Frazier and Rayner (1990), Frazier (1999)





Ambiguity Packing in Biology and Language

• two states are "temporarily equivalent" if their rightmost unpaired brackets are the same



psycholinguistic evidence (eye-tracking experiments):

delayed disambiguation

John and Mary had 2 papers John and Mary had 2 papers

Frazier and Rayner (1990), Frazier (1999)





Ambiguity Packing in Biology and Language

• two states are "temporarily equivalent" if their rightmost unpaired brackets are the same



psycholinguistic evidence (eye-tracking experiments):

delayed disambiguation

John and Mary had 2 papers each John and Mary had 2 papers together

Frazier and Rayner (1990), Frazier (1999)





Our Linear-Time Prediction is Much Faster...





with even slightly better prediction accuracy!!





... and Also More Accurate!

esp. on longer sequences and long-range base pairs







... and Also More Accurate!



Example: B. Subtilis 16S rRNA (length: 1,552nt)

World's Fastest RNA Structure Prediction Server

LinearFold Web Server (beta)				
- Interactive demo				
Add a sequence				
Paste or type your sequence here:				
CCUUCGAUAGCUCAGCUGGUAGAGCGGAGGACUGUAGAUCCUUAGGUCGCUGGUUCGAUUCCGGCUCGAAGGACCA				
Samples ▼ Or upload a file in FASTA format: Choose File No file chosen				
Set beam size				
Beam size (1-200): 100				
 Choose model(s) LinearFold-C (using <u>CONTRAfold v2.0</u> machine-learned model, Do et al 2006) LinearFold-V (using <u>Vienna RNAfold</u> thermodynamic model, Lorenz et al 2011, with parameters from Mathews et al 2004) 				
Run >> Reset				

http://linearfold.eecs.oregonstate.edu:8080/





Incremental Parsing <=> Incremental Folding

- incrementally parsable



Fast Structure Prediction Enables RNA Design



detecting active TB using RNA design

which needs our fast RNA folding

RNA sequence

GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA







Professor Rhiju Das Stanford Medical School



EteRNA game (RNA design)







tRNA 2-dimensional structure

tRNA 3-dimensional structure







Other Work, Recap, and Vision

Other Work, Recap, and Vision

























scene parsing



50

Other Work: Incremental Semantic Parsing

- parse NL (e.g., English) into a formal meaning representation type-driven parsing (formal semantics) + polymorphism (PL theory) • future work: (simultaneously) translating NL into PL (e.g., SQL)

		What is the capital of the largest state	by area	1?
step	action	stack	queue	typing
0	-	ϕ	what	
1–3	skip	$ \phi$	capital	
4	sh _{capital}	capital :st→ct	of	
7	sh _{largest}	capital :st \rightarrow ct argmax : ('a \rightarrow t) \rightarrow ('a \rightarrow i) \rightarrow 'a	state	
8	sh _{state}	capital:st \rightarrow ct argmax:('a \rightarrow t) \rightarrow ('a \rightarrow i) \rightarrow 'a state:st \rightarrow t	by	
9	re _へ	capital:st \rightarrow ct (argmax state):(st \rightarrow i) \rightarrow st	by	binding: 'a = st
11	sh _{area}	capital:st \rightarrow ct (argmax state):(st \rightarrow i) \rightarrow st size:lo \rightarrow i	?	
12	rea	capital:st→ct (argmax state size):st	?	st <: lo \Rightarrow
				$ (lo \rightarrow i) <: (st \rightarrow i)$
13	re _へ	(capital (argmax state size)):ct	?	

f the largest state	by area?
---------------------	----------





linear-time search algorithms

structured prediction with deep learning





Longer-Term Vision

efficiently analyze and generate sequences with hierarchical structures: natural language, RNA/proteins, programming languages, music, etc.

> grammar formalisms (context-free & beyond)

protein folding

self.plural is an lambda function with an argument *n*, which returns result of boolean expression n not equal to I



= lambda n: int(n!=1) self.plural

NL <=> PL translation



5-Year Vision in Natural Language Processing









5-Year Vision in Natural Language Processing

- simultaneous translation: from speech-to-text to speech-to-speech
 - Incremental text-to-speech synthesis (language production is also incremental!)
 - incremental predictive parsing on the source side (improve reordering)
 - incremental predictive parsing on the target side (improve prosody)









5-Year Vision in Natural Language Processing

- simultaneous translation: from speech-to-text to speech-to-speech
 - Incremental text-to-speech synthesis (language production is also incremental!)
 - incremental predictive parsing on the source side (improve reordering)
 - incremental predictive parsing on the target side (improve prosody)
- helping the language-impaired with incremental parsing and simultaneous MT
 - simultaneous English <=> ASL translation; intelligent input system










5-Year Vision in Natural Language Processing

- simultaneous translation: from speech-to-text to speech-to-speech
 - Incremental text-to-speech synthesis (language production is also incremental!)
 - incremental predictive parsing on the source side (improve reordering)
 - incremental predictive parsing on the target side (improve prosody)
- helping the language-impaired with incremental parsing and simultaneous MT
 - simultaneous English <=> ASL translation; intelligent input system
- online grammatical error correction; automatic or computer-aided writing











5-Year Vision in Natural Language Processing

- simultaneous translation: from speech-to-text to speech-to-speech
 - Incremental text-to-speech synthesis (language production is also incremental!)
 - incremental predictive parsing on the source side (improve reordering)
 - incremental predictive parsing on the target side (improve prosody)
- helping the language-impaired with incremental parsing and simultaneous MT
 - simultaneous English <=> ASL translation; intelligent input system
- online grammatical error correction; automatic or computer-aided writing
- incremental semantic parsing & code generation: NL=>PL (e.g. SQL)
 - also the inverse problem of PL = > NL translation (comment generation)











5-Year Vision in Natural Language Processing

- simultaneous translation: from speech-to-text to speech-to-speech
 - Incremental text-to-speech synthesis (language production is also incremental!)
 - incremental predictive parsing on the source side (improve reordering)
 - incremental predictive parsing on the target side (improve prosody)
- simultaneous English <=> ASL translation; intelligent input system
- online grammatical error correction; automatic or computer-aided writing
- incremental semantic parsing & code generation: NL=>PL (e.g. SQL)
 - also the inverse problem of PL = > NL translation (comment generation)
- how do incremental predictive parsers compare with psycholinguistics data? 53















5-Year Vision in

- simultaneous translation: from speech-to-text to speech-to-speech
 - Incremental text-to-speech synthesis (language production is also incremental!)
 - incremental predictive parsing on the source side (improve reordering)
 - incremental predictive parsing on the target side (improve prosody)
- simultaneous English <=> ASL translation; intelligent input system
- online grammatical error correction; automatic or computer-aided writing
- incremental semantic parsing & code generation: NL=>PL (e.g. SQL)
 - also the inverse problem of PL = > NL translation (comment generation)
- how do incremental predictive parsers compare with psycholinguistics data? 53

helping people communicate across linguistic and accessibility barriers















5-Year Vision in Computational Biology

- Inear-time incremental folding: from RNA to protein structure prediction
- predicting crossing structures in RNAs/proteins: use linear-time parsing and *mildly* context-sensitive grammars (polynomial-time parsable)
- how does our beam search compare to real incremental folding in nature?



Chomsky Hierarchy



reestablishing the forgotten link between 5-Year Vision in computational linguistics and structural biology

- Inear-time incremental folding: from RNA to protein structure prediction
- predicting crossing structures in RNAs/proteins: use linear-time parsing and *mildly* context-sensitive grammars (polynomial-time parsable)
- how does our beam search compare to real incremental folding in nature?











非常感谢您 来 听 我 的 演讲

Thank you very much for listening to my speech

Bush met Putin in Moscow

source language sentence



natural language sentence



eat sushi with tuna from Japan geggaauageuegguagageaegaegaeguegegguegegguegegguegegguegegguegegguegegguegegguegegguegegueg

RNA sequence





非常感谢您 来 听 我 的 演讲



Thank you very much for listening to my speech

Happy Chinese New Year!

target-language sequence

布什在莫斯科与普京会晤

Bush met Putin in Moscow

source language sentence





natural language sentence



secondary structure

RNA sequence







Backup Slide



Chinese input:	
Pinyin:	
Word-by-Word Translation:	
Simultaneous Translation (wait-3):	
Simultaneous Tranlation (wait-5):	
Simultaneous Translation (wait-3 + revise):	







Backup Slide



Chinese input:	
Pinyin:	
Word-by-Word Translation:	
Simultaneous Translation (wait-3):	
Simultaneous Tranlation (wait-5):	
Simultaneous Translation (wait-3 + revise):	







Translation with Noisy Speech Input

neural MT is fragile, and automatic speech recognition output is noisy

Clean Input Output of Transformer	目前已发现 <u>有</u> 109人 at present, 109 peopl
Noisy Input	目前已发现 <u>又</u> 109人
Output of Transformer	the hpv has been fou
Output of Our Method	so far, 109 people ha

by contrast, is very robust to homophone noises thanks to phonetic information.

- our work (Liu et al, ArXiv 2018): Robust Neural MT using phonetic information
 - 、死亡,另有57人获救
 - le have been found dead and 57 have been rescued
 - 、死亡,另有57人获救
 - and dead so far and 57 have been saved
 - we been found dead and 57 others have been rescued
- Table 1: The translation results on Mandarin sentences without and with homophone noises. The word '有' (yǒu, "have") in clean input is replaced by one of its homophone, ' χ ' (you, "again"), to form a noisy input. This seemingly minor change completely fools the Transformer to generate something irrelvant ("hpv"). Our method,





