

Green AI

Roy Schwartz

Hebrew University of Jerusalem

CS Colloquium, EPFL
December 2022



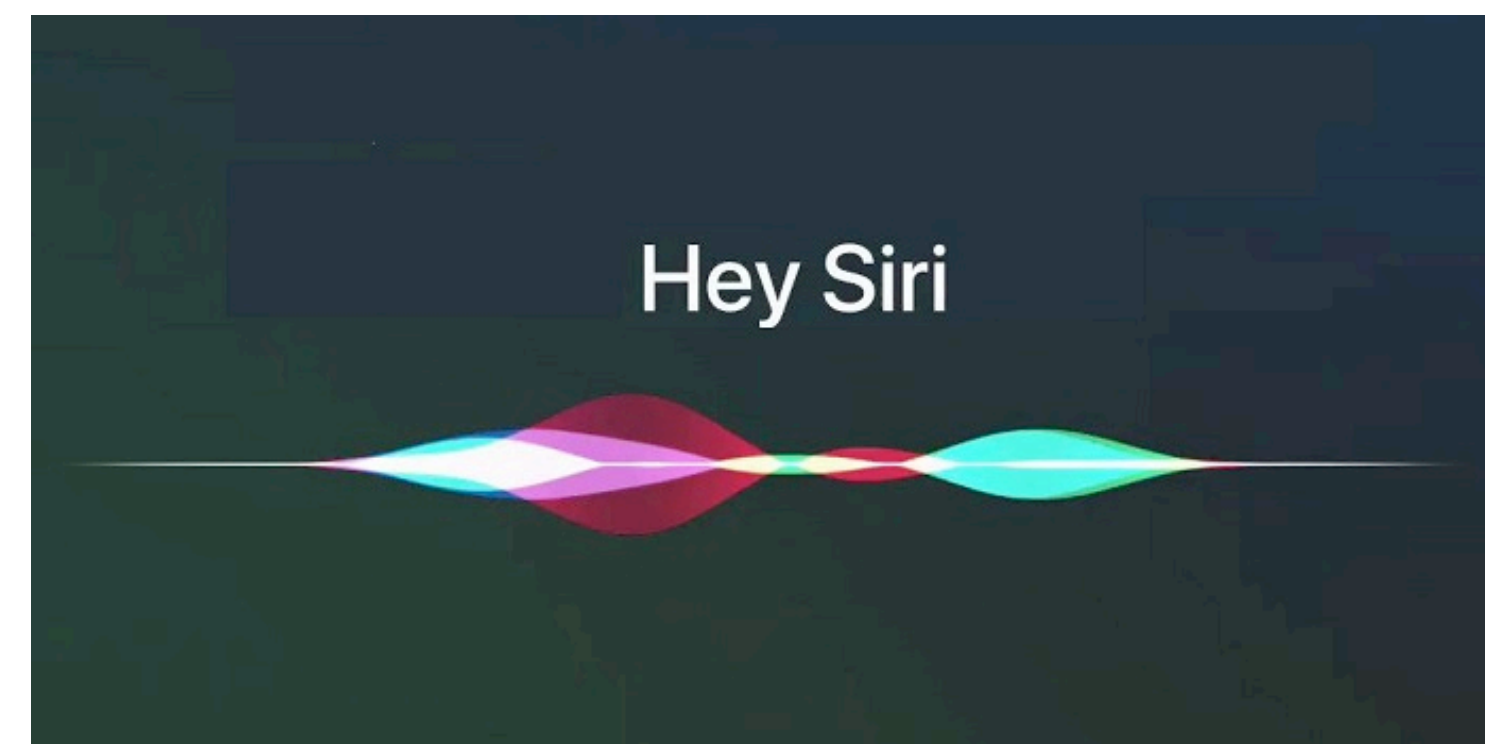
THE HEBREW
UNIVERSITY
OF JERUSALEM



AI Today

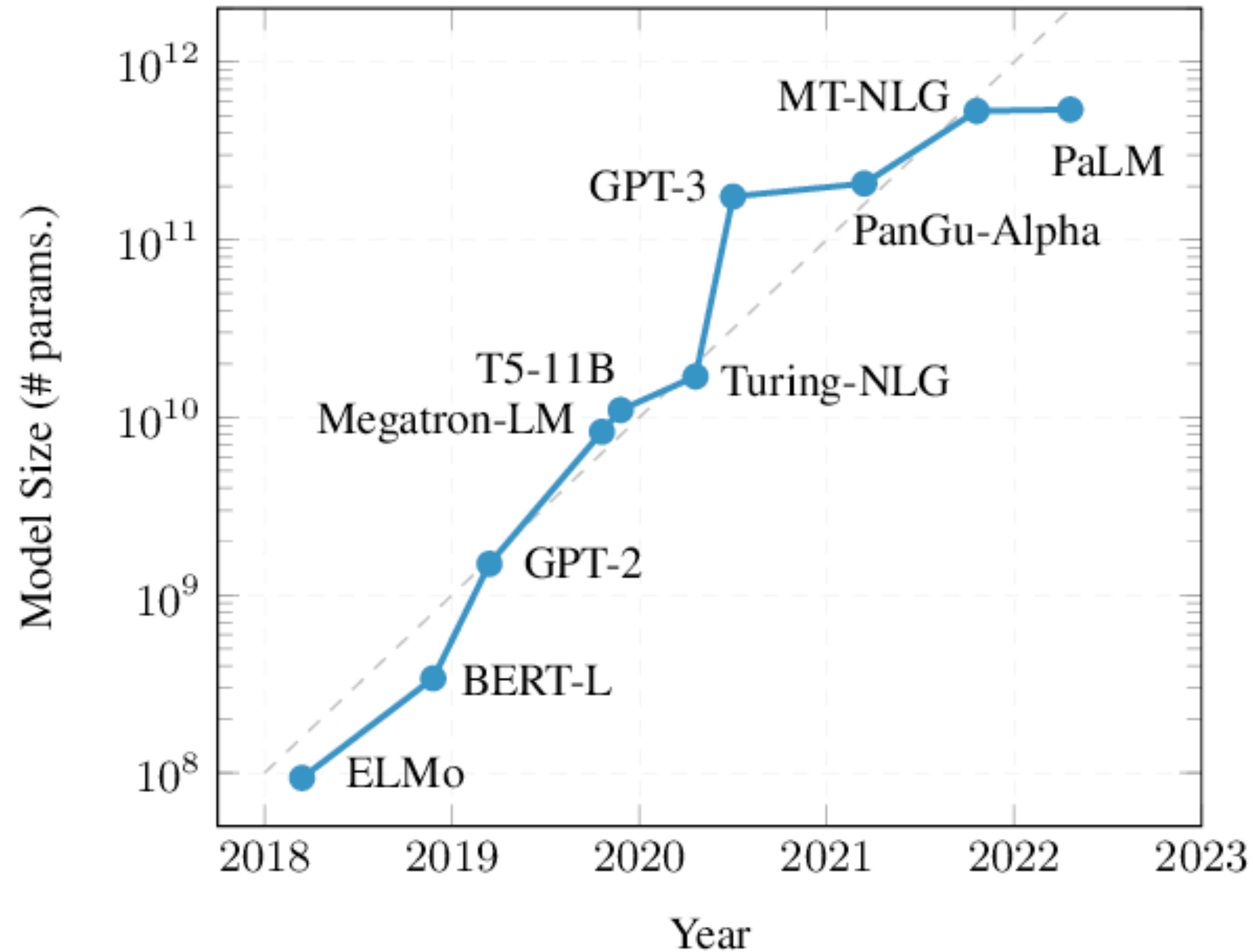


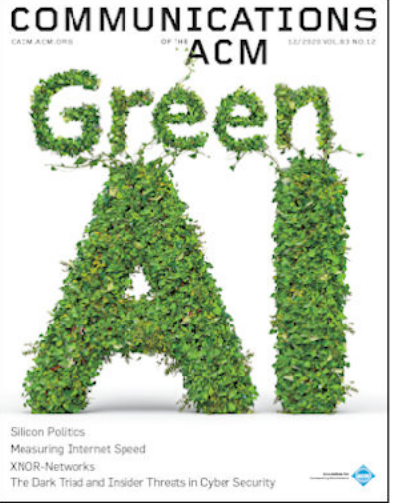
ChatGPT		
☀ Examples	⚡ Capabilities	⚠ Limitations
"Explain quantum computing in simple terms"	Remembers what user said earlier in the conversation	May occasionally generate incorrect information
"Got any creative ideas for a 10 year old's birthday?"	Allows user to provide follow-up corrections	May occasionally produce harmful instructions or biased content
"How do I make an HTTP request in Javascript?"	Trained to decline inappropriate requests	Limited knowledge of world and events after 2021



Scaling

5,000X in 4 Years

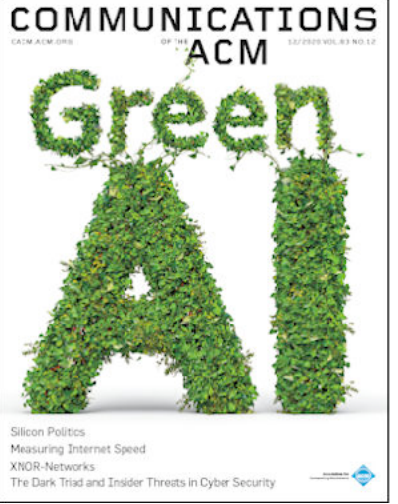




Green AI

Schwartz*, Dodge*, Smith & Etzioni, CACM 2020

- **Red AI**
- Problems: inclusiveness, environment



Green AI

Schwartz*, Dodge*, Smith & Etzioni, CACM 2020

- **Red AI**

- Problems: inclusiveness, environment

- **Green AI**

- Enhance **reporting** of computational budgets

- Add a *price-tag* for scientific results

- Promote **efficiency** as a core evaluation for AI

- **In addition to** accuracy

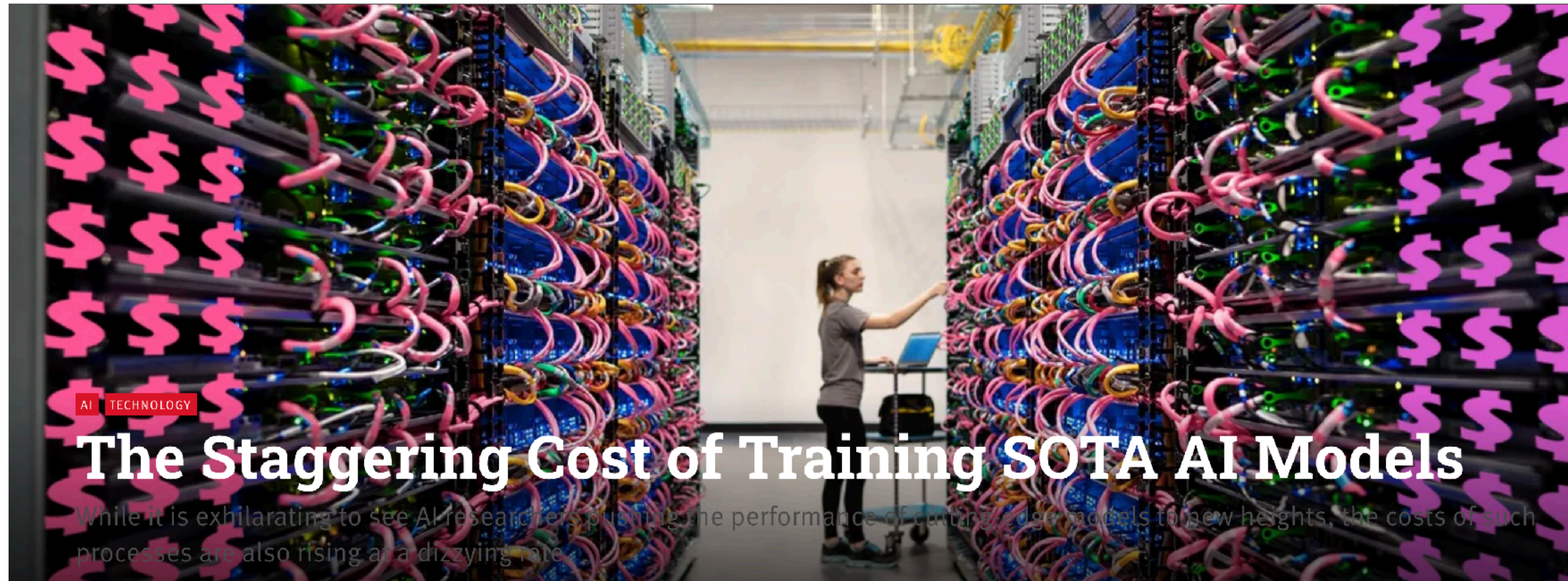


Problems with **Scaling** Inclusiveness

Synced

AI TECHNOLOGY & INDUSTRY REVIEW

FEATURE ▾ INDUSTRY ▾ TECHNOLOGY COMMUNITY ▾ ABOUT US ▾ REPORT CONTRIBUTE TO SYNCED REVIEW



<https://syncedreview.com/2019/06/27/the-staggering-cost-of-training-sota-ai-models/>

Training Costs

- BERT (Devlin et al, 2019) was trained on **16** Cloud TPUs for **4** days
- RoBERTa (Liu et al., 2019) was trained on **1024** V100 GPUs for approximately **1** day
- PaLM (Chowdhery et al., 2022) was trained on **6144** TPU v4 chips for **50** days and **3072** TPU v4 chips for **15** days

Training Costs

- BERT (Devlin et al, 2019) was trained on **16** Cloud TPUs for **4** days
- RoBERTa (Liu et al., 2019) was trained on **1024** V100 GPUs for approximately **1** day
- PaLM (Chowdhery et al., 2022) was trained on **6144** TPU v4 chips for **50** days and **3072** TPU v4 chips for **15** days

We need better reporting!



Number of Authors

Number of Authors

Language Models are Few-Shot Learners

Tom B. Brown* Benjamin Mann* Nick Ryder* Melanie Subbiah*
Jared Kaplan† Prafulla Dhariwal Arvind Neelakantan Pranav Shyam Girish Sastry
Amanda Askell Sandhini Agarwal Ariel Herbert-Voss Gretchen Krueger Tom Henighan
Rewon Child Aditya Ramesh Daniel M. Ziegler Jeffrey Wu Clemens Winter
Christopher Hesse Mark Chen Eric Sigler Mateusz Litwin Scott Gray
Benjamin Chess Jack Clark Christopher Berner
Sam McCandlish Alec Radford Ilya Sutskever Dario Amodei

OpenAI

Number of Authors

Language Models are Few-Shot Learners

Tom B. Brown* Benjamin Mann* Nick Ryder* Melanie Subbiah*
Jared Kaplan† Prafulla Dhariwal Arvind Neelakantan Pranav Shyam Girish Sastry
Amanda Askell Sandhini Agarwal Ariel Herbert-Voss Gretchen Krueger Tom Henighan
Rewon Child Aditya Ramesh Daniel M. Ziegler Jeffrey Wu Clemens Winter
Christopher Hesse Mark Chen Eric Sigler Mateusz Litwin Scott Gray
Benjamin Chess Jack Clark Christopher Berner
Sam McCandlish Alec Radford Ilya Sutskever Dario Amodei

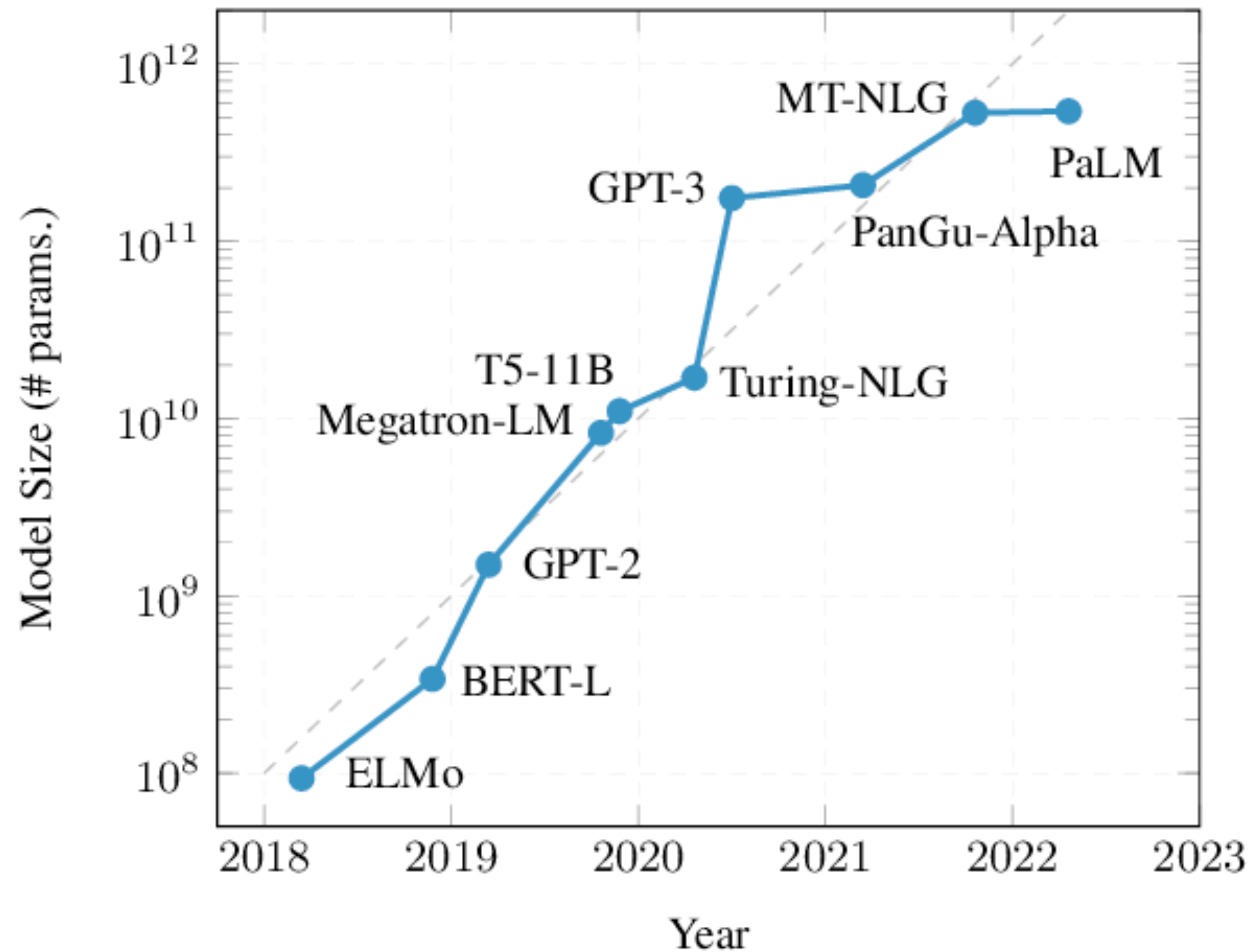
OpenAI

PaLM: Scaling Language Modeling with Pathways

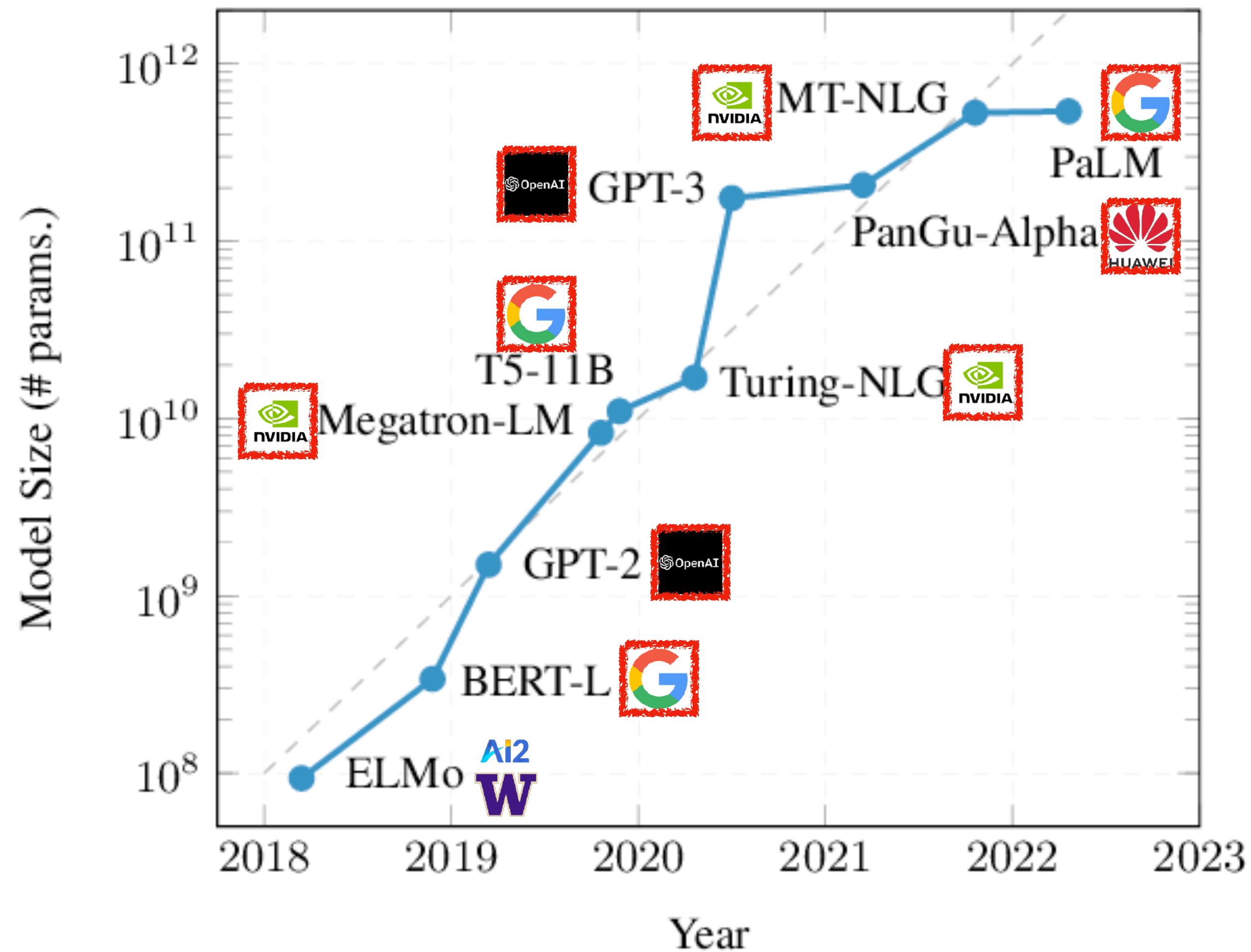
Aakanksha Chowdhery* Sharan Narang* Jacob Devlin*
Maarten Bosma Gaurav Mishra Adam Roberts Paul Barham
Hyung Won Chung Charles Sutton Sebastian Gehrmann Parker Schuh Kensen Shi
Sasha Tsvyashchenko Joshua Maynez Abhishek Rao† Parker Barnes Yi Tay
Noam Shazeer† Vinodkumar Prabhakaran Emily Reif Nan Du Ben Hutchinson
Reiner Pope James Bradbury Jacob Austin Michael Isard Guy Gur-Ari
Pengcheng Yin Toju Duke Anselm Levskaya Sanjay Ghemawat Sunipa Dev
Henryk Michalewski Xavier Garcia Vedant Misra Kevin Robinson Liam Fedus
Denny Zhou Daphne Ippolito David Luan† Hyeontaek Lim Barret Zoph
Alexander Spiridonov Ryan Sepassi David Dohan Shivani Agrawal Mark Omernick
Andrew M. Dai Thanumalayan Sankaranarayanan Pillai Marie Pellat Aitor Lewkowycz
Erica Moreira Rewon Child Oleksandr Polozov† Katherine Lee Zongwei Zhou
Xuezhi Wang Brennan Saeta Mark Diaz Orhan Firat Michele Catasta† Jason Wei
Kathy Meier-Hellstern Douglas Eck Jeff Dean Slav Petrov Noah Fiedel

Google Research

It's a Rich Man's World



It's a Rich Man's World



Problems with **Scaling** Environment

Consumption	CO₂e (lbs)
Air travel, 1 person, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000
Training one model (GPU)	
NLP pipeline (parsing, SRL)	39
w/ tuning & experiments	78,468
Transformer (big)	192
w/ neural arch. search	626,155

Strubell et al. (2019)

Problems with **Scaling** Environment

Consumption	CO₂e (lbs)
Air travel, 1 person, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000
Training one model (GPU)	
NLP pipeline (parsing, SRL)	39
w/ tuning & experiments	78,468
Transformer (big)	192
w/ neural arch. search	626,155

Strubell et al. (2019)

*Is AI really creating an
environmental problem?*

Google's Answer: No!

BLOG ›

Good News About the Carbon Footprint of Machine Learning Training

TUESDAY, FEBRUARY 15, 2022

Posted by David Patterson, Distinguished Engineer, Google Research, Brain Team

*Strubell et al.'s energy estimate for NAS ended up **18.7X** too high for the average organization (...) and **88X** off in emissions for energy-efficient organizations like Google*

Google's Answer: No!

BLOG ›

Good News About the Carbon Footprint of Machine Learning Training

TUESDAY, FEBRUARY 15, 2022

Posted by David Patterson, Distinguished Engineer, Google Research, Brain Team

*Strubell et al.'s energy estimate for NAS ended up **18.7X** too high for the average organization (...) and **88X** off in emissions for energy-efficient organizations like Google*

We need better reporting!



Our Answer: Maybe?

Measuring the Carbon Intensity of AI in Cloud Instances

JESSE DODGE, Allen Institute for AI, USA

TAYLOR PREWITT, University of Washington, USA

REMI TACHET DES COMBES, Microsoft Research Montreal, USA

ERIKA ODMARK, Microsoft, USA

ROY SCHWARTZ, Hebrew University of Jerusalem, Israel

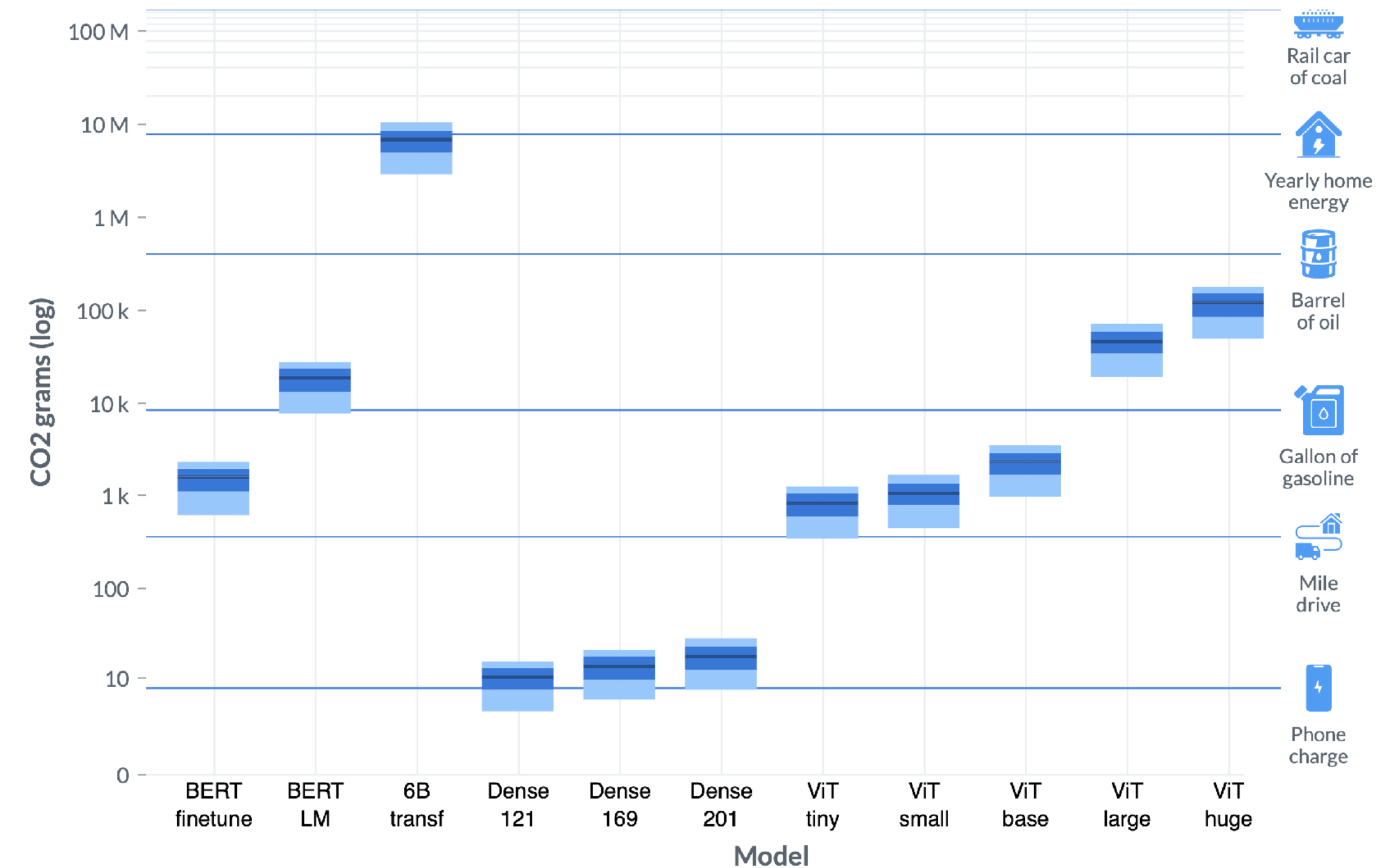
EMMA STRUBELL, Carnegie Mellon University, USA

ALEXANDRA SASHA LUCCIONI, Hugging Face, USA

NOAH A. SMITH, Allen Institute for AI and University of Washington, USA

NICOLE DECARIO, Allen Institute for AI, USA

WILL BUCHANAN, Microsoft, USA

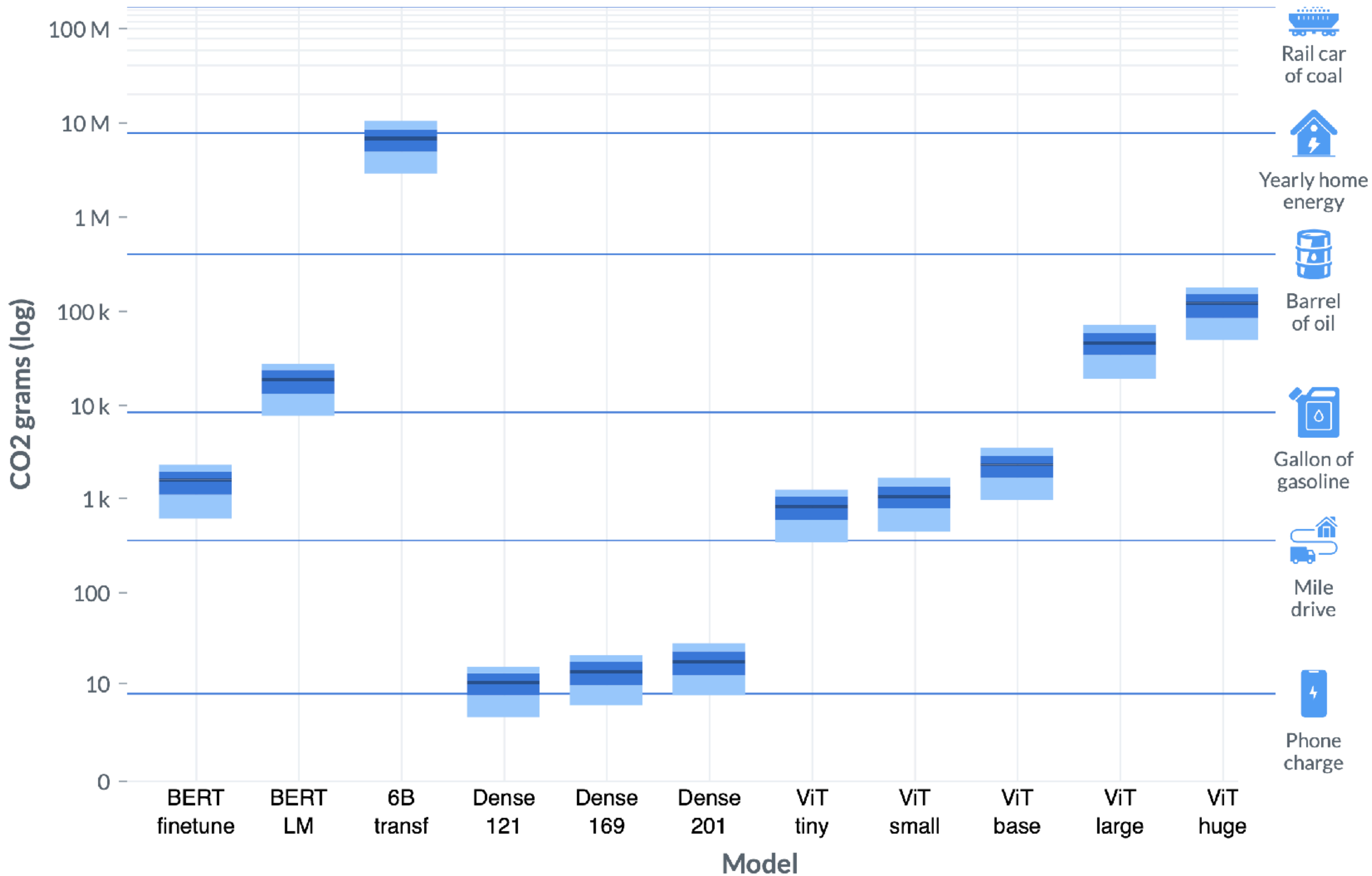


CO2 Relative Size Comparison



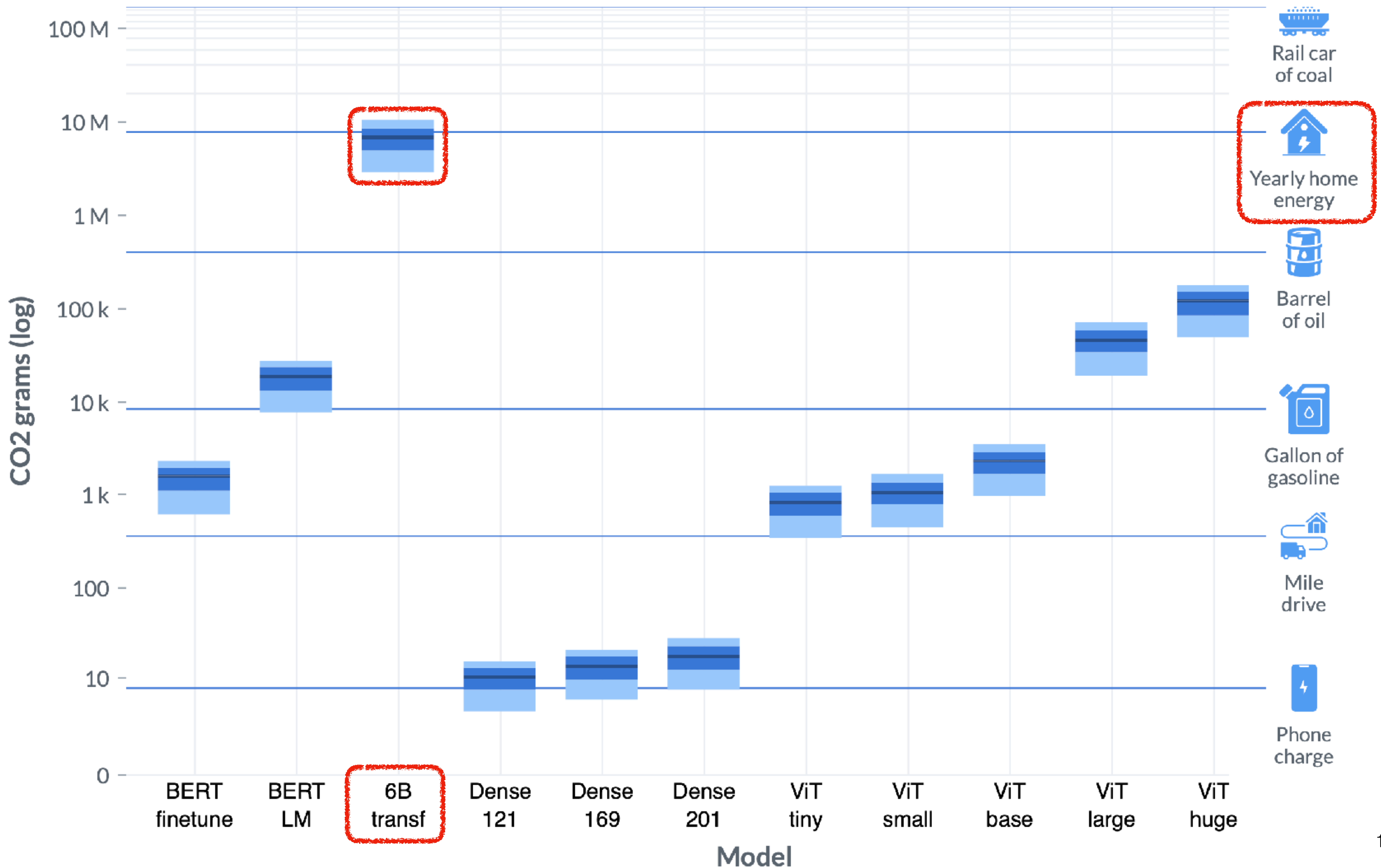
Measur

- JESSE I
- TAYLO
- REMI T
- ERIKA
- ROY SC
- EMMA
- ALEXA
- NOAH
- NICOL
- WILL B



Meast

JESSE I
TAYLO
REMI T
ERIKA
ROY SC
EMMA
ALEXA
NOAH
NICOL
WILL B



AI and the Environment



- Evidence around the **most expensive experiments**
 - More recent models consume 2-3 orders of magnitude more CO₂ (Luccioni et al., 2022)
 - But these are typically run very few times

AI and the Environment



- Evidence around the **most expensive experiments**
 - More recent models consume 2-3 orders of magnitude more CO₂ (Luccioni et al., 2022)
 - But these are typically run very few times
- What about “normal” experiments?
 - Much **cheaper**, but run **hundreds / thousands of times a day**?

AI and the Environment



- Evidence around the **most expensive experiments**
 - More recent models consume 2-3 orders of magnitude more CO₂ (Luccioni et al., 2022)
 - But these are typically run very few times
- What about “normal” experiments?
 - Much **cheaper**, but run **hundreds / thousands of times a day**?
- What about **inference** operations?
 - Very **cheap** (though increasingly more expensive)
 - Run **billions of times a day**?
 - 80-90% of AI computation is spent on inference

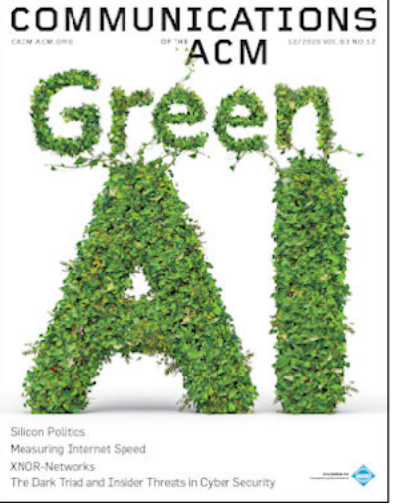
AI and the Environment



- Evidence around the **most expensive experiments**
 - More recent models consume 2-3 orders of magnitude more CO₂ (Luccioni et al., 2022)
 - But these are typically run very few times
- What about “normal” experiments?
 - Much **cheaper**, but run **hundreds / thousands of times a day**?
- What about **inference** operations?
 - Very **cheap** (though increasingly more expensive)
 - Run **billions of times a day**?
 - 80-90% of AI computation is spent on inference

We need better reporting!





Green AI

Schwartz*, Dodge*, Smith & Etzioni, CACM 2020

- **Red AI**

- Problems: inclusiveness, environment

- **Green AI**

- Enhance **reporting** of computational budgets

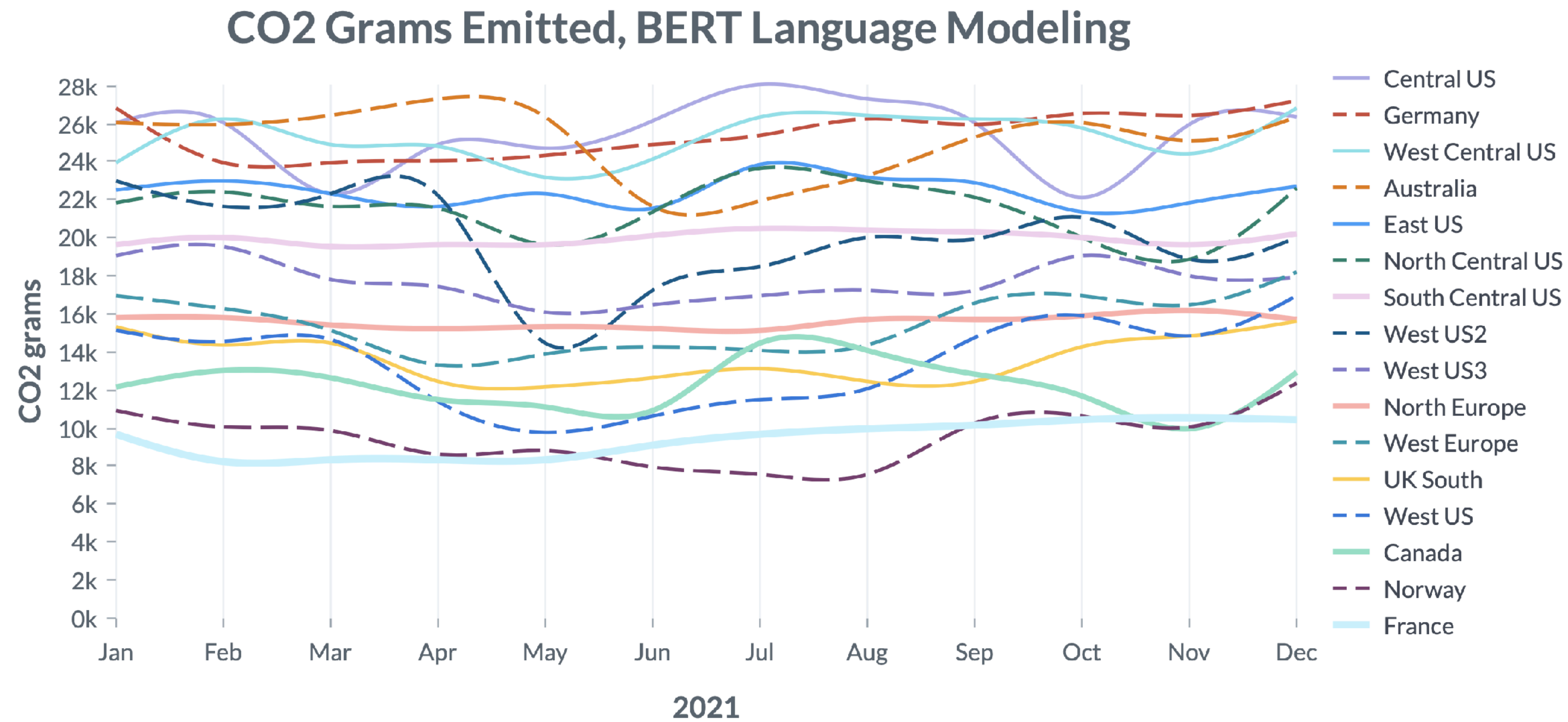
- Add a *price-tag* for scientific results

- Promote **efficiency** as a core evaluation for AI

- **In addition to** accuracy

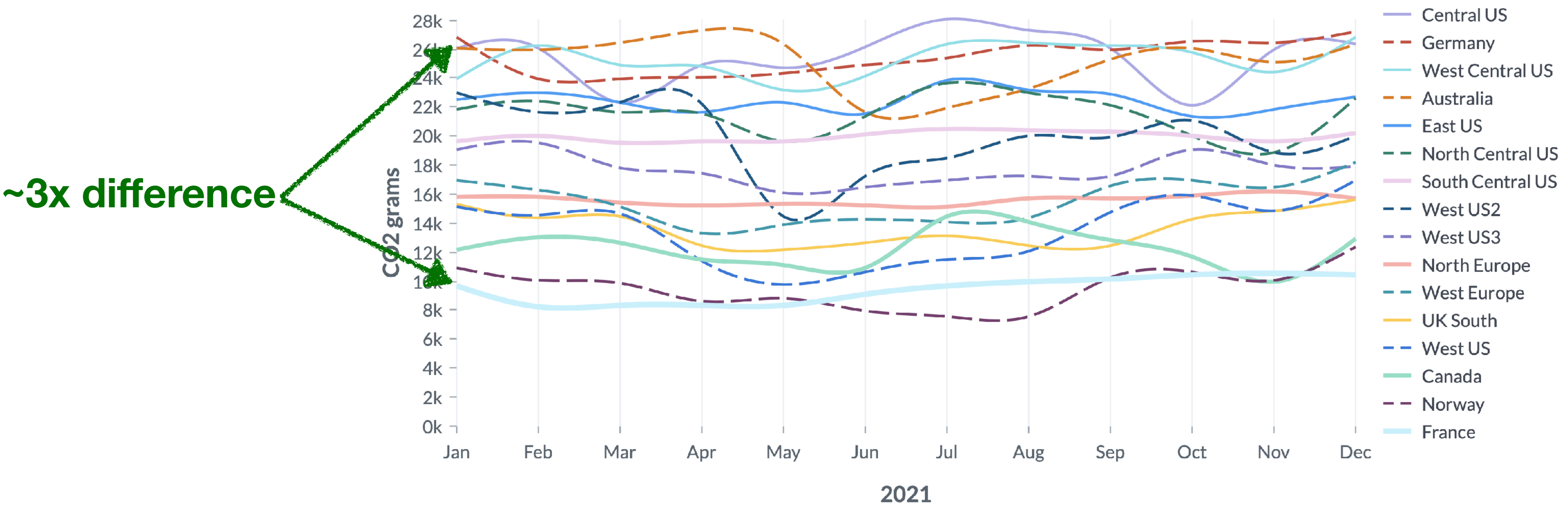


Cloud Location Matters



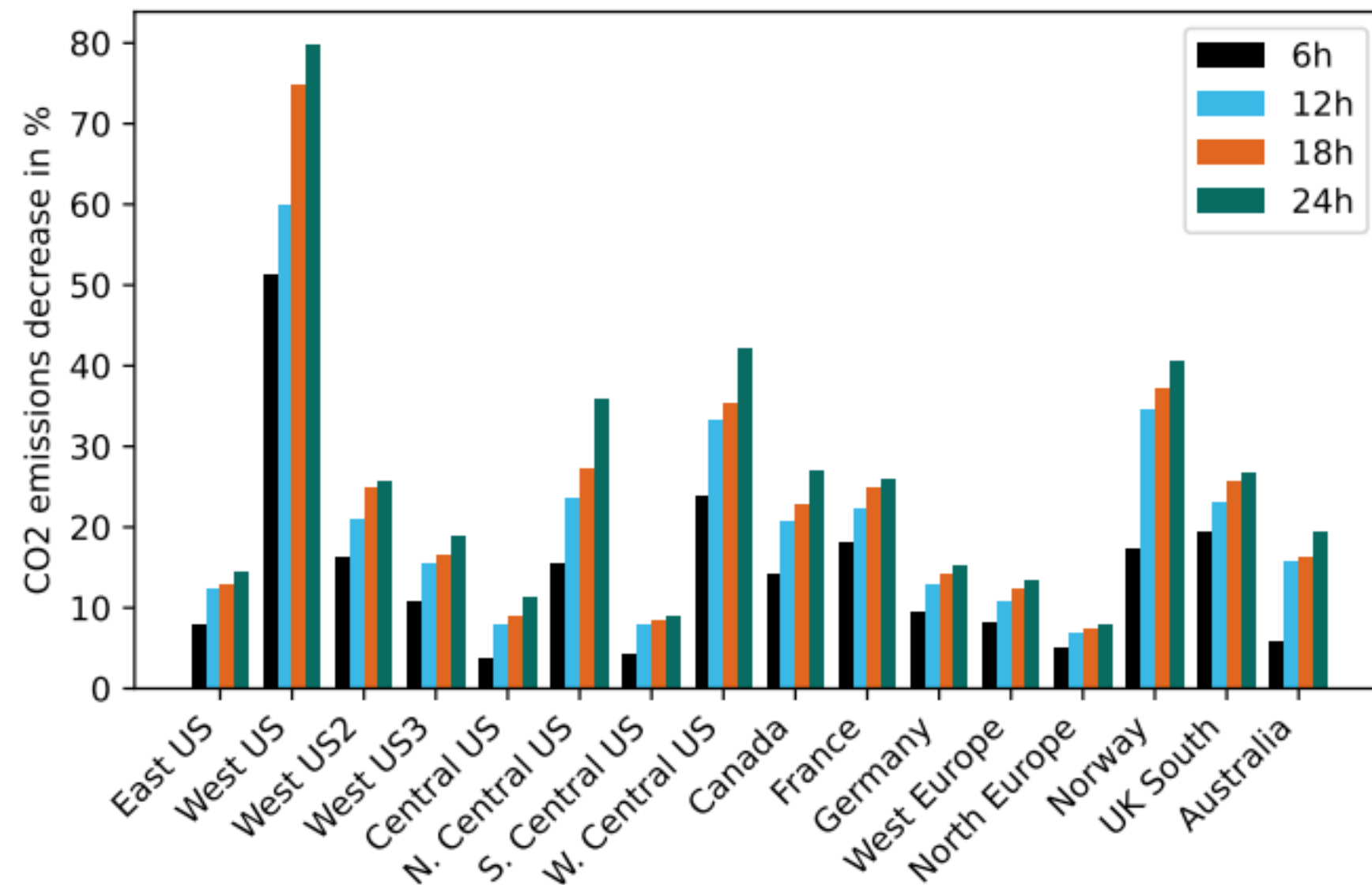
Cloud Location Matters

CO2 Grams Emitted, BERT Language Modeling

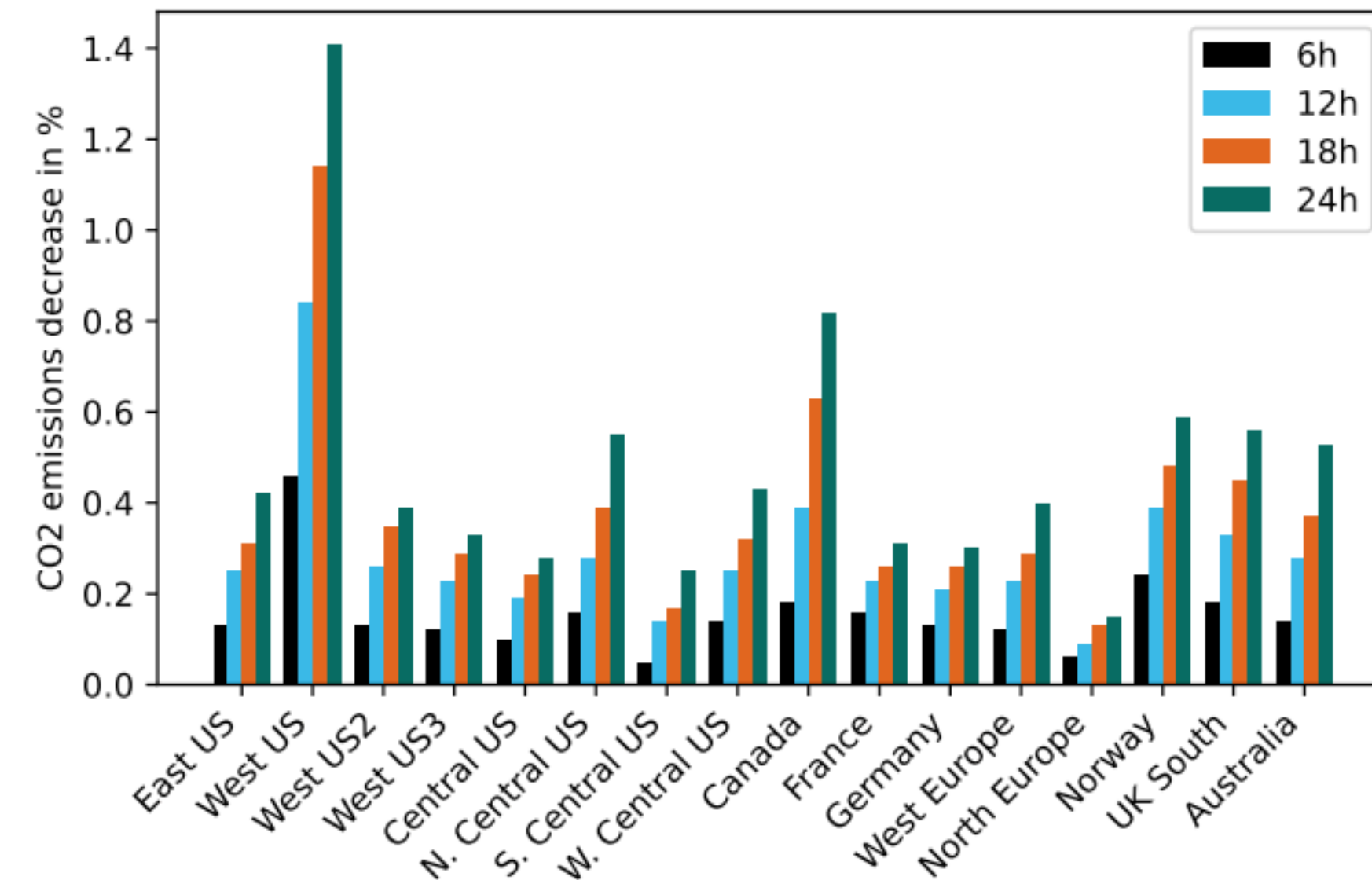


Time of Day Matters

Potential Saving with *Flexible Start*



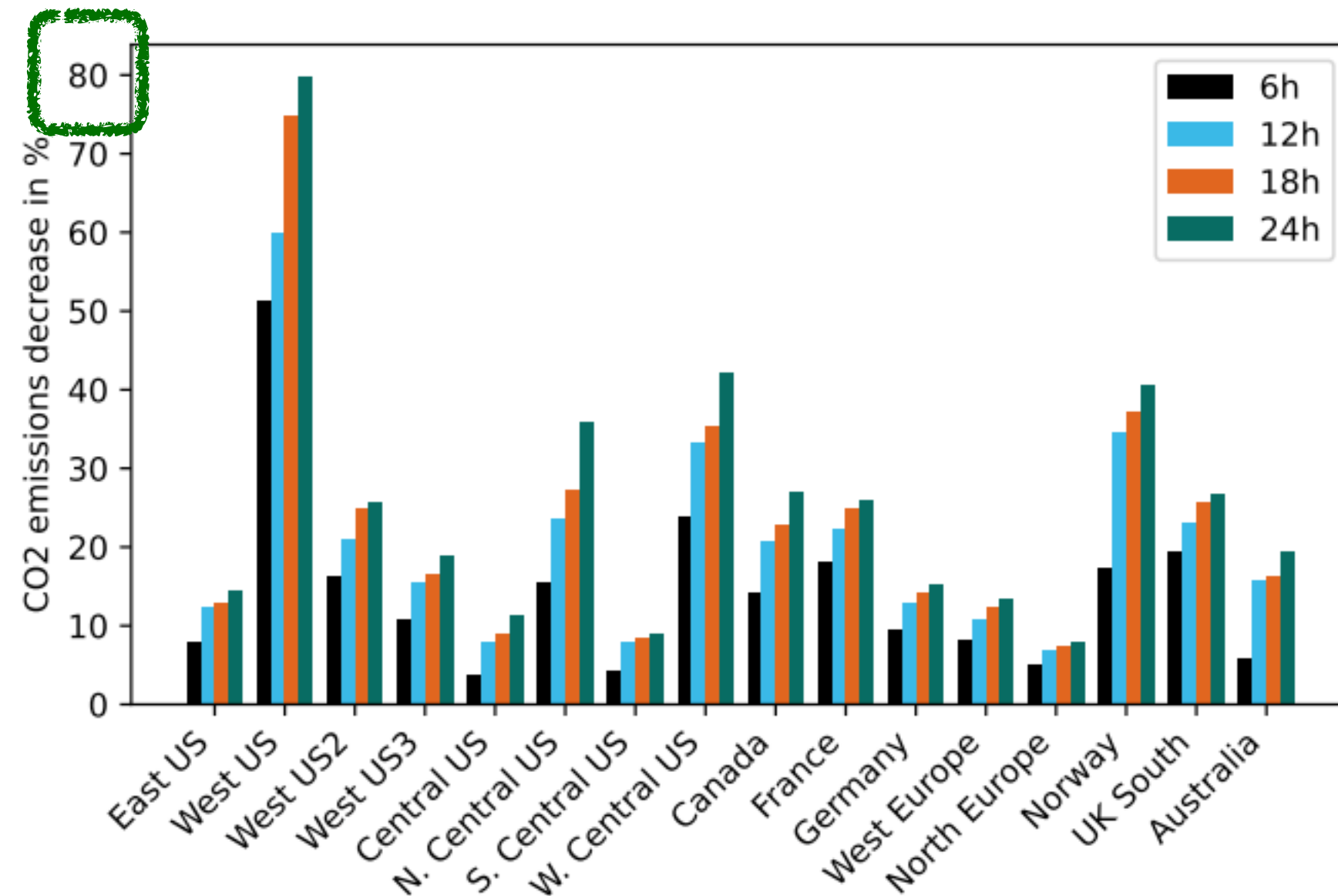
(a) *Flexible Start* optimization for Dense 201.



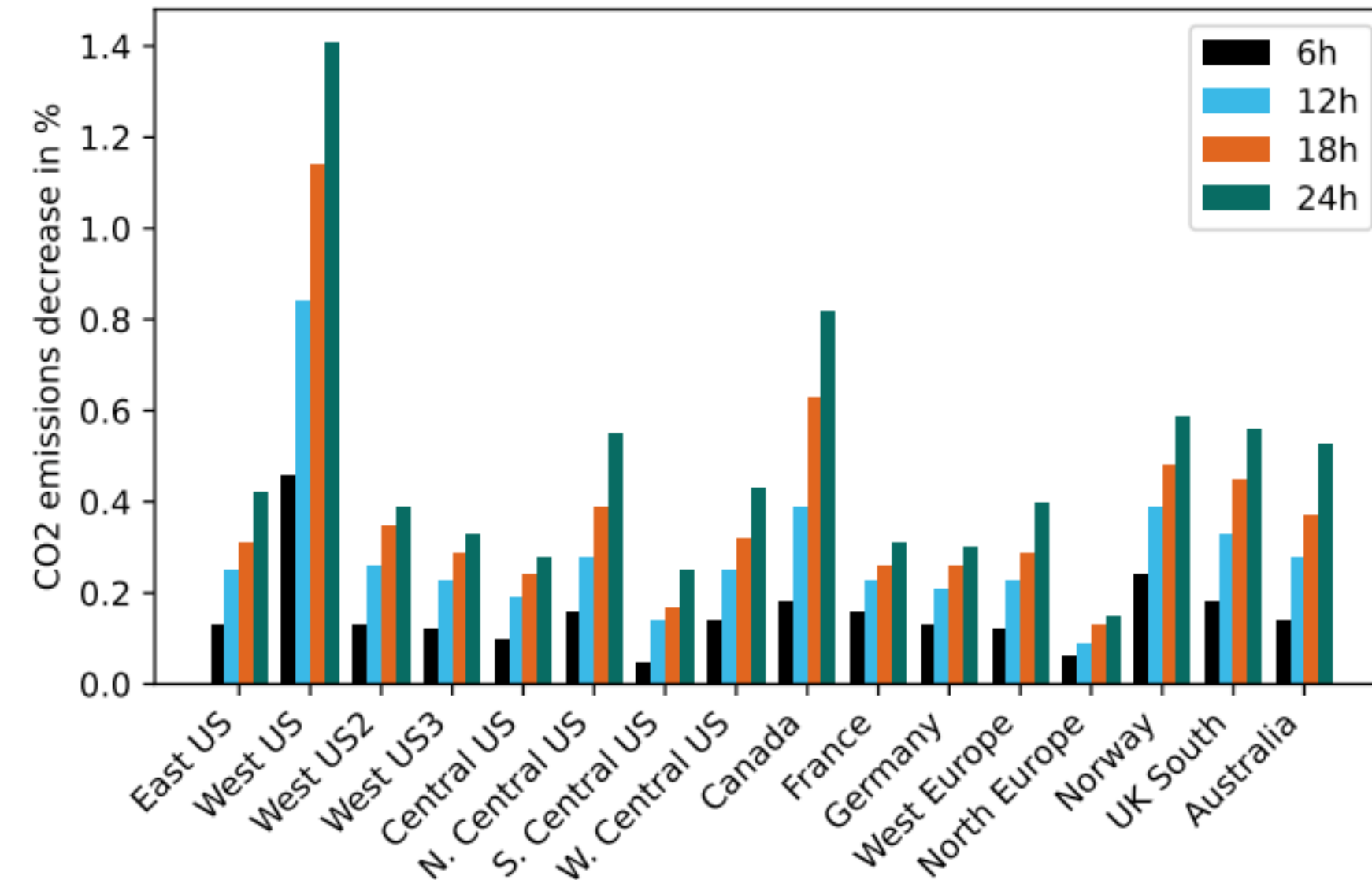
(b) *Flexible Start* optimization for 6B parameters Transformer.

Time of Day Matters

Potential Saving with *Flexible Start*



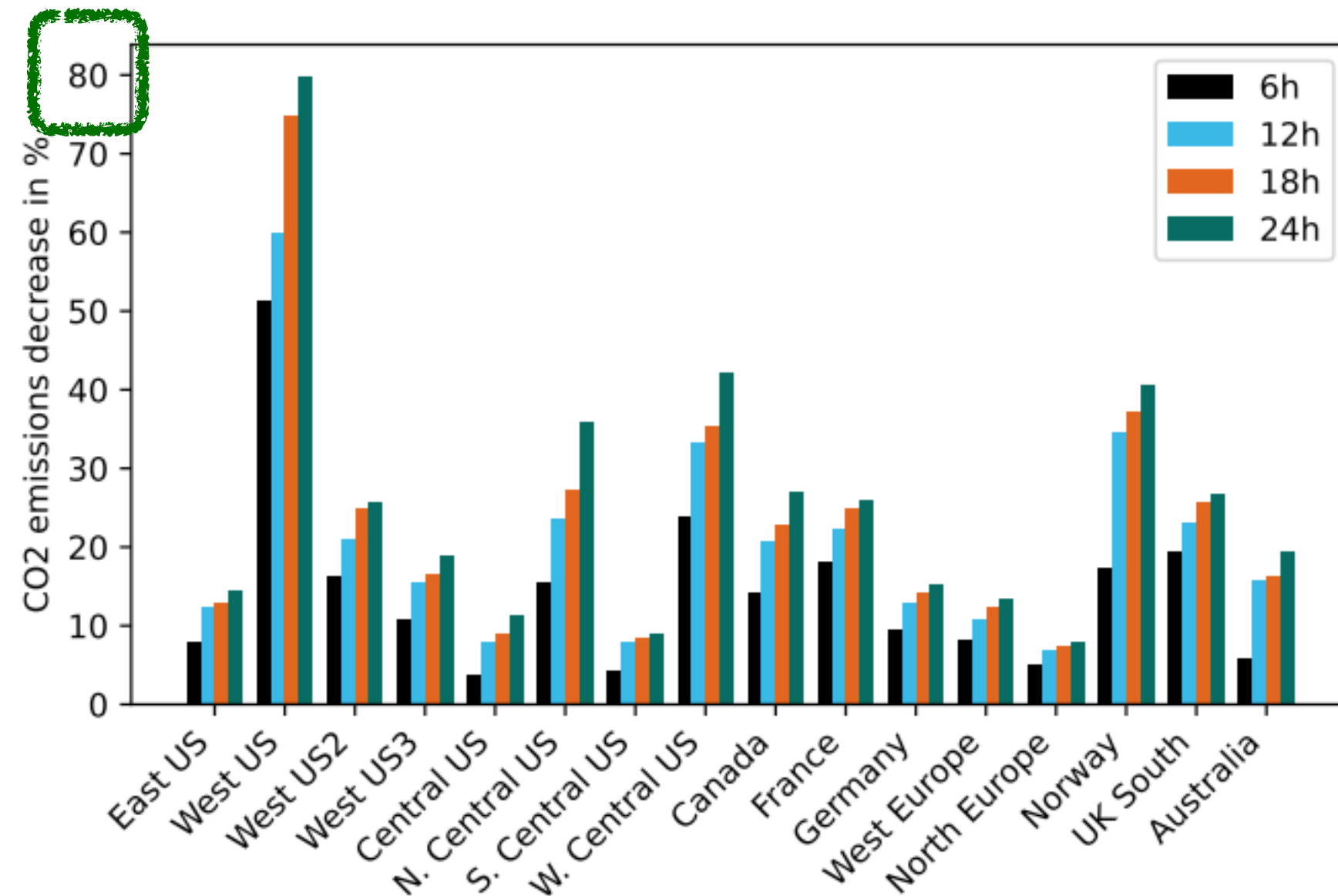
(a) *Flexible Start* optimization for Dense 201.



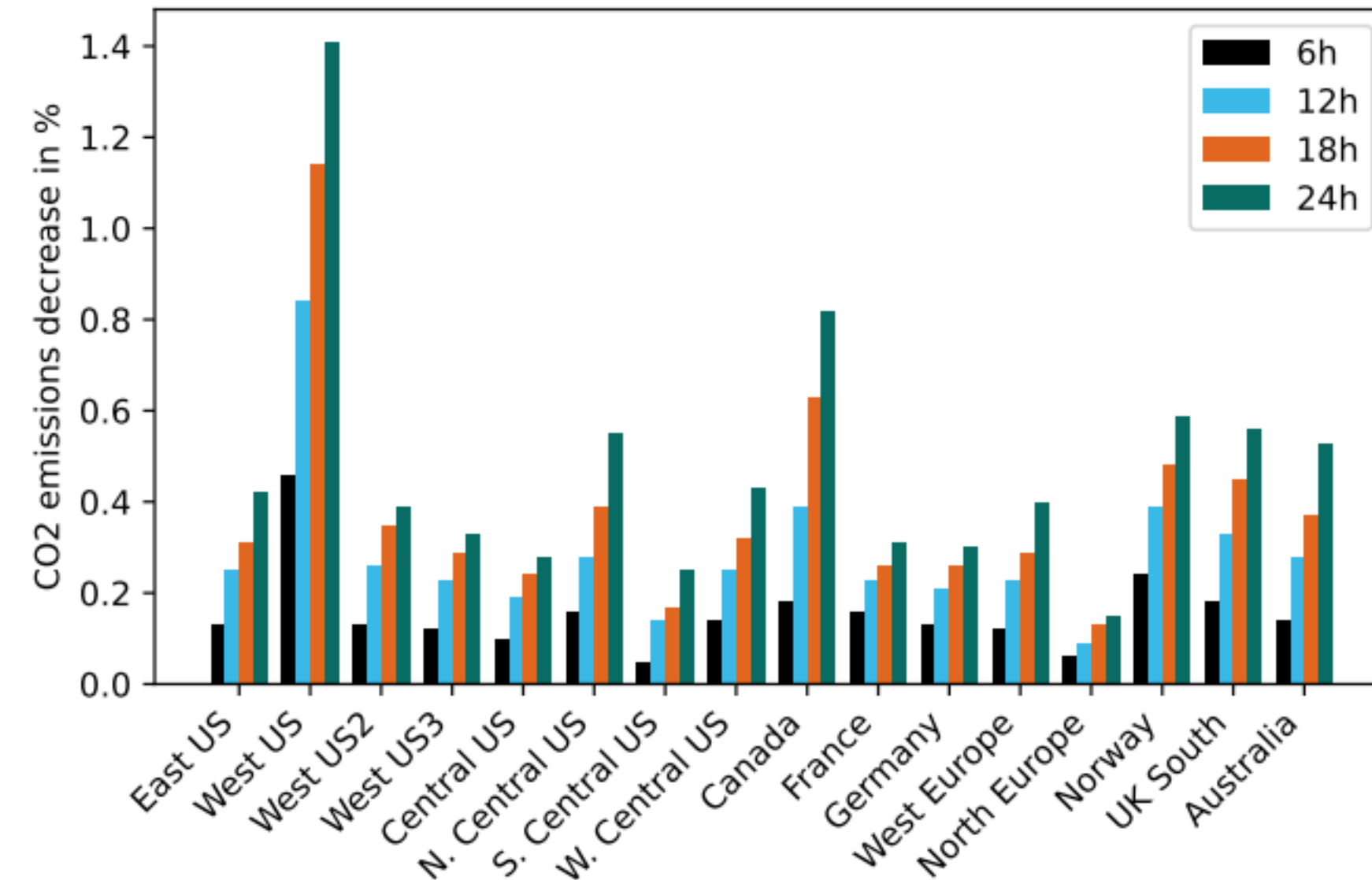
(b) *Flexible Start* optimization for 6B parameters Transformer.

Time of Day Matters

Potential Saving with *Flexible Start*



(a) *Flexible Start* optimization for Dense 201.



(b) *Flexible Start* optimization for 6B parameters Transformer.

We need better reporting!



Stop training large models?

Large Models are Important

- Push the limits of SOTA
- Released large pre-trained models **save compute**
- Large models are potentially faster to train
 - Li et al. (2020)

Large Models are Important

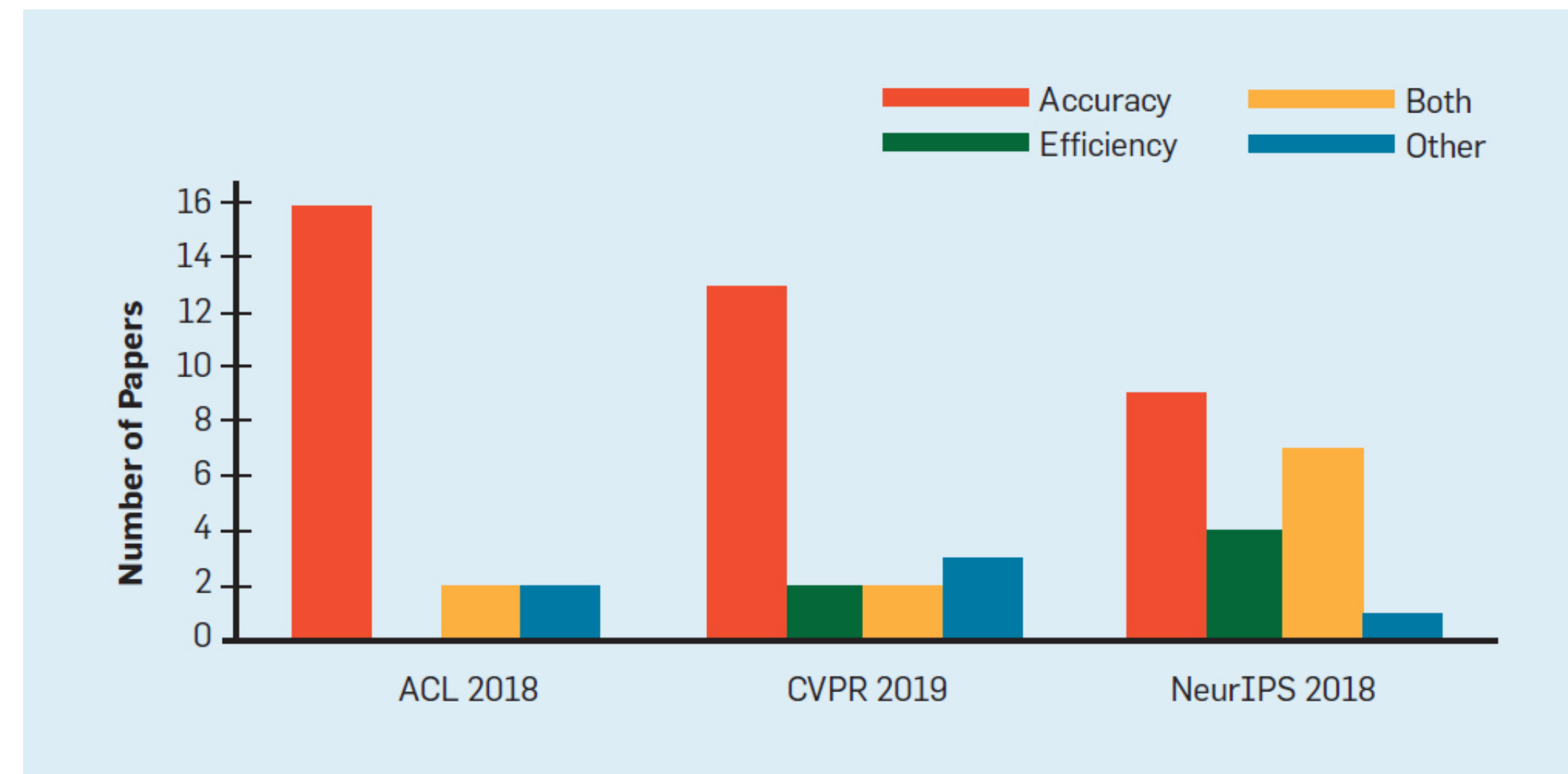
- Push the limits of SOTA
- Released large pre-trained models **save compute**
- Large models are potentially faster to train
 - Li et al. (2020)
- But, **large models** have **concerning side effects**
 - Inclusiveness, environment

Large Models are Important

- Push the limits of SOTA
- Released large pre-trained models **save compute**
- Large models are potentially faster to train
 - Li et al. (2020)
- But, **large models** have **concerning side effects**
 - Inclusiveness, environment
- Our goal is to **mitigate these side effects**

Accuracy or Efficiency?

Accuracy or Efficiency?



S. et al. (2020)



Efficient NLP

Yuki Arase, Osaka University, arase@ist.osaka-u.ac.jp

Phil Blunsom, Oxford University, phil.blunsom@cs.ox.ac.uk

Mona Diab, George Washington University & Facebook AI, mtdiab@gwu.edu

Jesse Dodge, Allen Institute for AI, jessed@allenai.org

Iryna Gurevych, Technical University of Darmstadt, gurevych@ukp.informatik.tu-darmstadt.de

Percy Liang, Stanford University, pliang@cs.stanford.edu

Colin Raffel, UNC & HuggingFace, craffel@gmail.com

Andreas Rücklé, Amazon, andreas@rueckle.net

Roy Schwartz, The Hebrew University of Jerusalem, roy.schwartz1@mail.huji.ac.il

Noah A. Smith, University of Washington & Allen Institute for AI, nasmith@cs.washington.edu

Emma Strubell, Carnegie Mellon University & Google, strubell@cmu.edu

Yue Zhang, Westlake University, yue.zhang@wias.org.cn



Efficient NLP

Yuki Arase, Osaka University, arase@ist.osaka-u.ac.jp

Phil Blunsom, Oxford University, phil.blunsom@cs.ox.ac.uk

Mona Diab, George Washington University & Facebook AI, mtdiab@gwu.edu

Jesse Dodge, Allen Institute for AI, jessed@allenai.org

Iryna Gurevych, Technical University of Darmstadt, gurevych@ukp.informatik.tu-darmstadt.de

Percy Liang, Stanford University, pliang@cs.stanford.edu

Colin Raffel, UNC & HuggingFace, craffel@gmail.com

Andreas Rücklé, Amazon, andreas@rueckle.net

Roy Schwartz, The Hebrew University of Jerusalem, roy.schwartz1@mail.huji.ac.il

Noah A. Smith, University of Washington & Allen Institute for AI, nasmith@cs.washington.edu

Emma Strubell, Carnegie Mellon University & Google, strubell@cmu.edu

Yue Zhang, Westlake University, yue.zhang@wias.org.cn

- Setting up conference areas that target **efficiency**



Efficient NLP

Yuki Arase, Osaka University, arase@ist.osaka-u.ac.jp

Phil Blunsom, Oxford University, phil.blunsom@cs.ox.ac.uk

Mona Diab, George Washington University & Facebook AI, mtdiab@gwu.edu

Jesse Dodge, Allen Institute for AI, jessed@allenai.org

Iryna Gurevych, Technical University of Darmstadt, gurevych@ukp.informatik.tu-darmstadt.de

Percy Liang, Stanford University, pliang@cs.stanford.edu

Colin Raffel, UNC & HuggingFace, craffel@gmail.com

Andreas Rücklé, Amazon, andreas@rueckle.net

Roy Schwartz, The Hebrew University of Jerusalem, roy.schwartz1@mail.huji.ac.il

Noah A. Smith, University of Washington & Allen Institute for AI, nasmith@cs.washington.edu

Emma Strubell, Carnegie Mellon University & Google, strubell@cmu.edu

Yue Zhang, Westlake University, yue.zhang@wias.org.cn

- Setting up conference areas that target **efficiency**
- Encouraging the **release** of trained models

You are here: **Program** » **Seminar Calendar** » Seminar Homepage

<https://www.dagstuhl.de/22232>

June 6 – 10 , 2022, Dagstuhl Seminar 22232

Efficient and Equitable Natural Language Processing in the Age of Deep Learning

Organizers

Jesse Dodge (AI2 – Seattle, US)

Iryna Gurevych (TU Darmstadt, DE)

Roy Schwartz (The Hebrew University of Jerusalem, IL)

Emma Strubell (Carnegie Mellon University – Pittsburgh, US)



Efficient Methods for Natural Language Processing: A Survey

**Marcos Treviso^{10*}, Tianchu Ji^{3*}, Ji-Ung Lee^{7*}, Betty van Aken⁸, Qingqing Cao²,
Manuel R. Ciosici⁹, Michael Hassid¹, Kenneth Heafield¹³, Sara Hooker⁵,
Pedro H. Martins¹⁰, André F. T. Martins¹⁰, Peter Milder³, Colin Raffel⁶,
Edwin Simpson⁴, Noam Slonim¹², Niranjan Balasubramanian³, Leon Derczynski¹¹, Roy Schwartz¹**

¹The Hebrew University of Jerusalem, ²University of Washington, ³Stony Brook University,

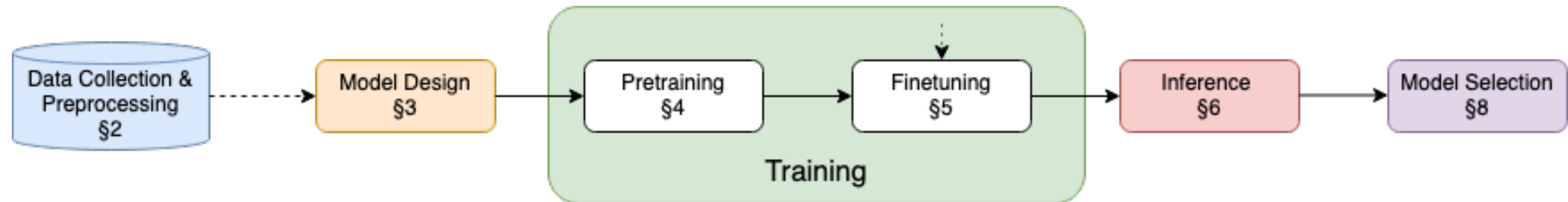
⁴University of Bristol, ⁵Cohere For AI, ⁶University of North Carolina at Chapel Hill,

⁷Technical University of Darmstadt, ⁸Berliner Hochschule für Technik,

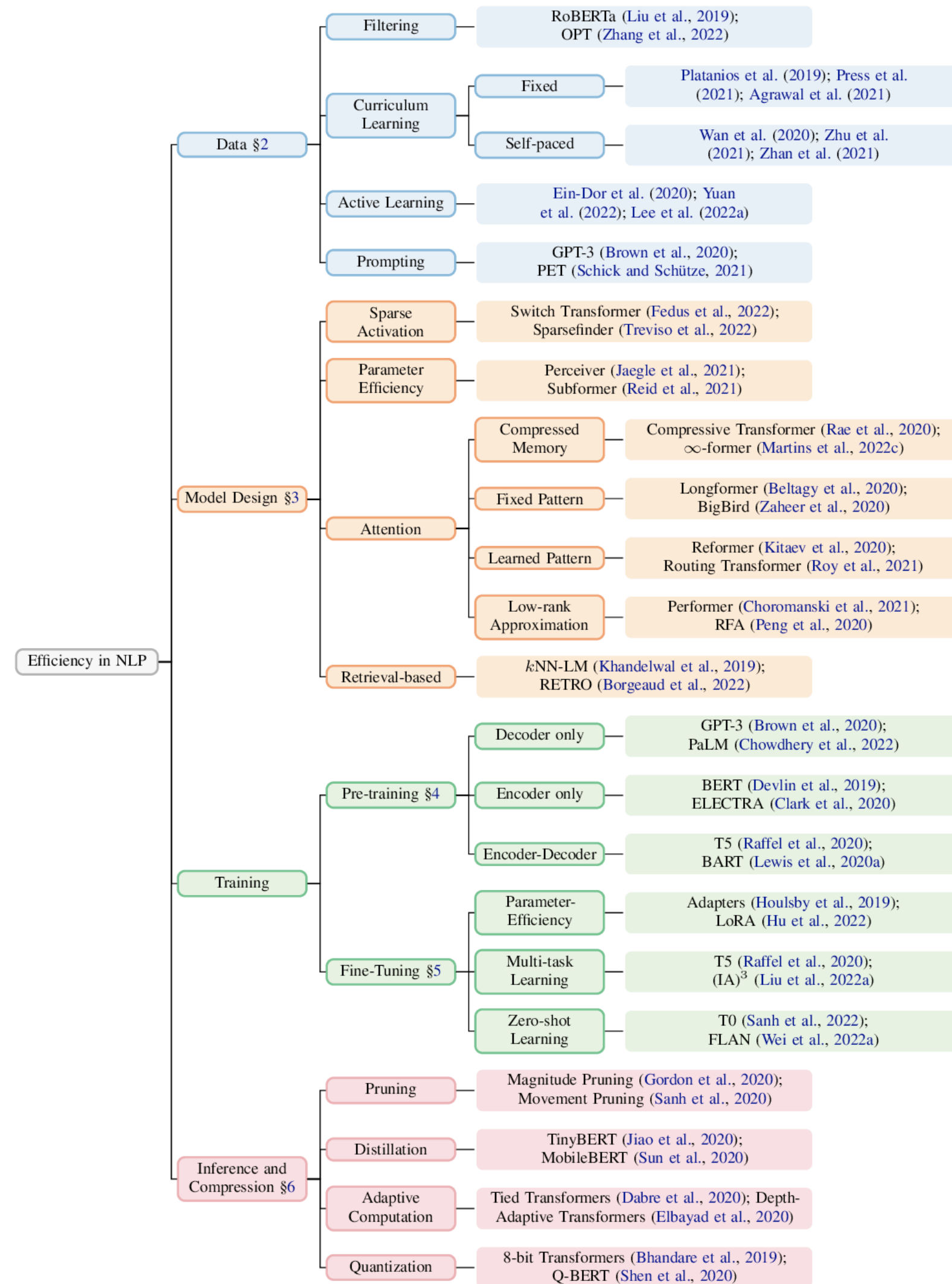
⁹University of Southern California, ¹⁰IST/University of Lisbon & Instituto de Telecomunicações,

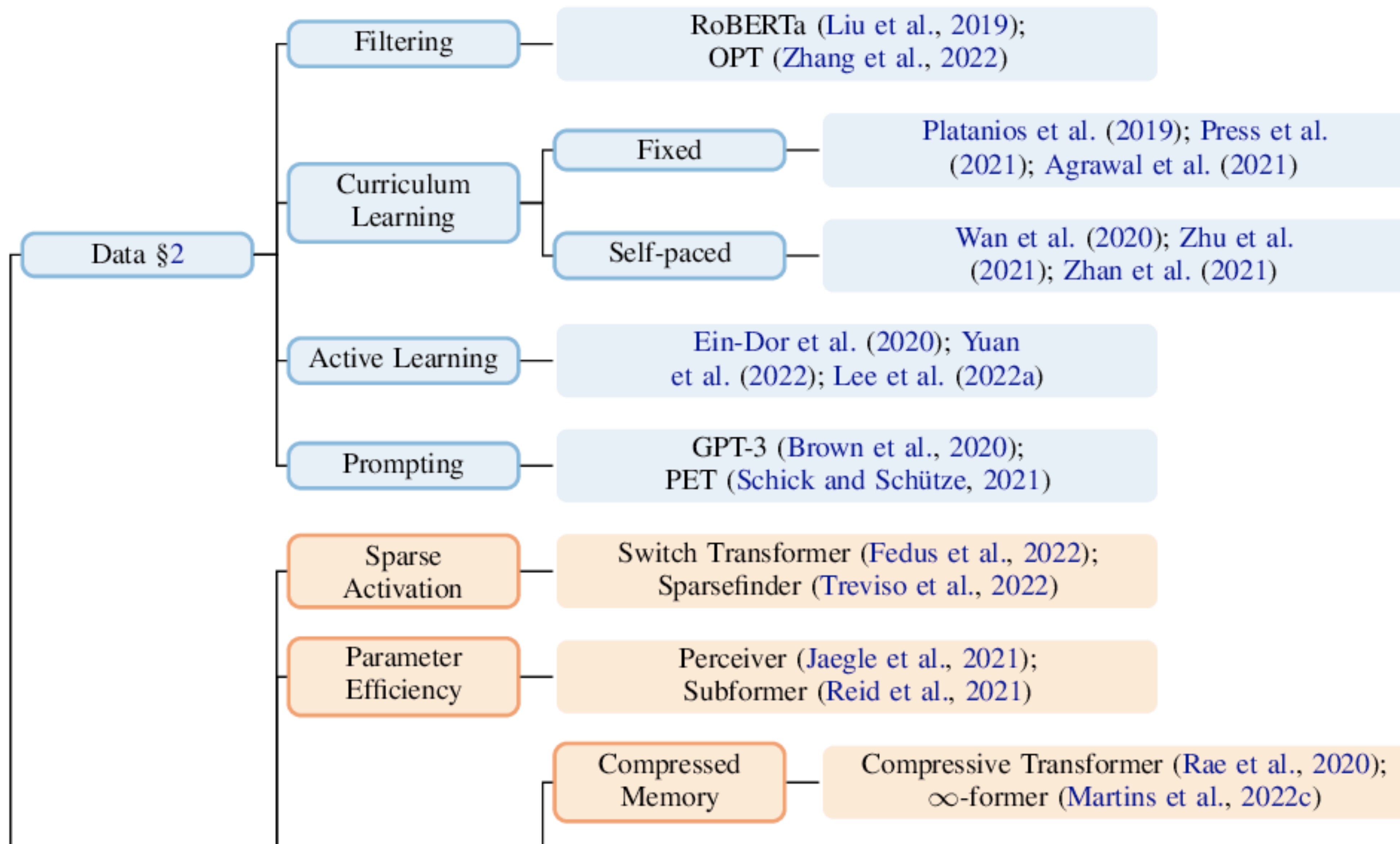
¹¹IT University of Copenhagen, ¹²IBM Research, ¹³University of Edinburgh

Efficient Methods in NLP



Efficient Methods in NLP





Filtering

- Non-text
 - Gibrish, HTML
- Text in other languages
- Foul text
- Typically done via simple, rule-based heuristics
 - Noisy process

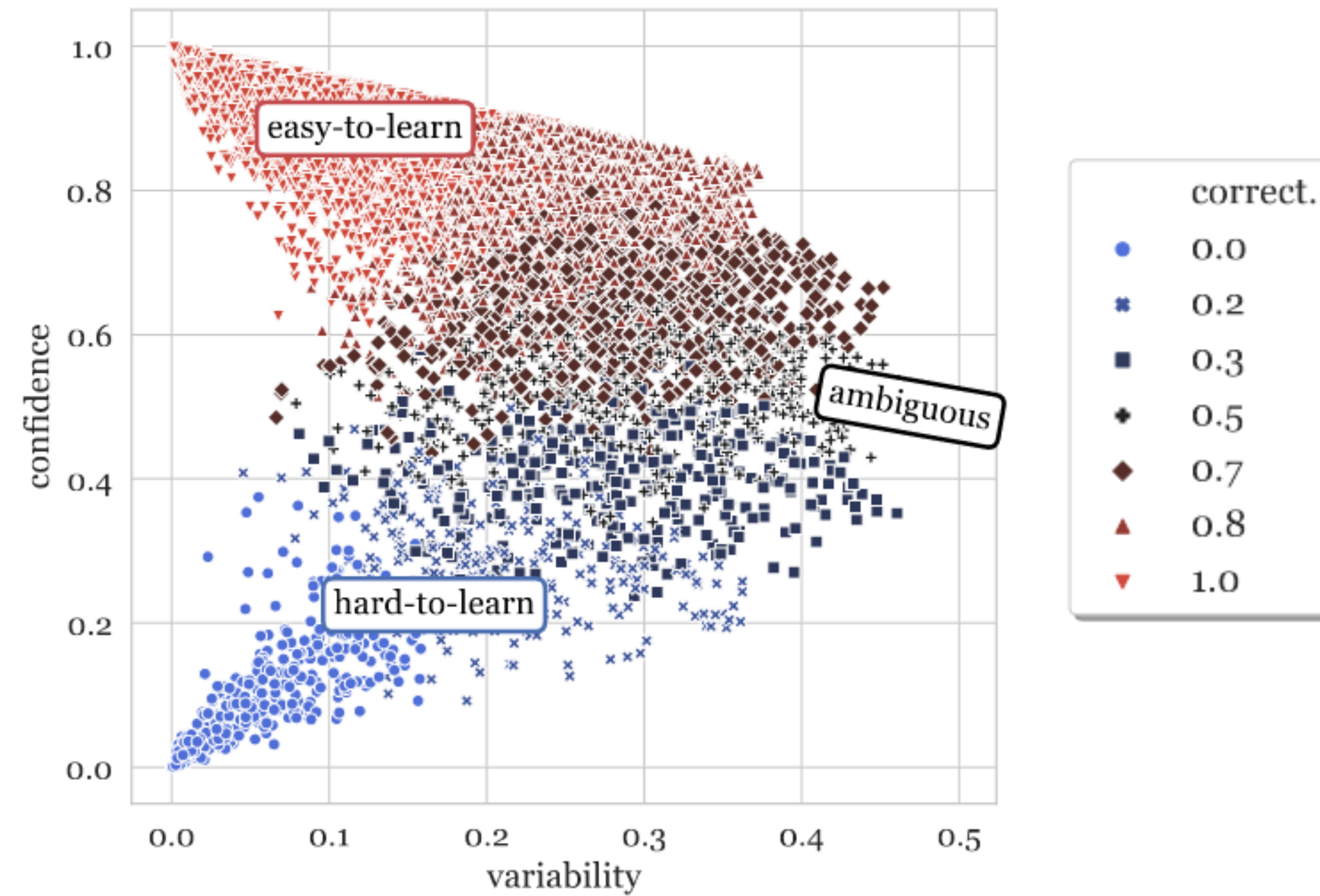
Smart Filtering

Swayamdipta, **S.** et al., EMNLP 2020

- Not all training instances contribute the same to learning
 - Some are “*easy-to-learn*”, others are more challenging



Dataset Map



Experiments

WinoGrande, RoBERTa-Large

	WINOg. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7 _{0.4}	87.6 _{0.6}

Experiments

WinoGrande, RoBERTa-Large

	WINOGr. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7 _{0.4}	87.6 _{0.6}

Experiments

WinoGrande, RoBERTa-Large

33% {

	WINOGr. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7 _{0.4}	87.6 _{0.6}

Experiments

WinoGrande, RoBERTa-Large

33%



	WINOg. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7 _{0.4}	87.6 _{0.6}

Experiments

WinoGrande, RoBERTa-Large

33% {

	WINOGr. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7_{0.4}	87.6_{0.6}

Experiments

WinoGrande, RoBERTa-Large

33%



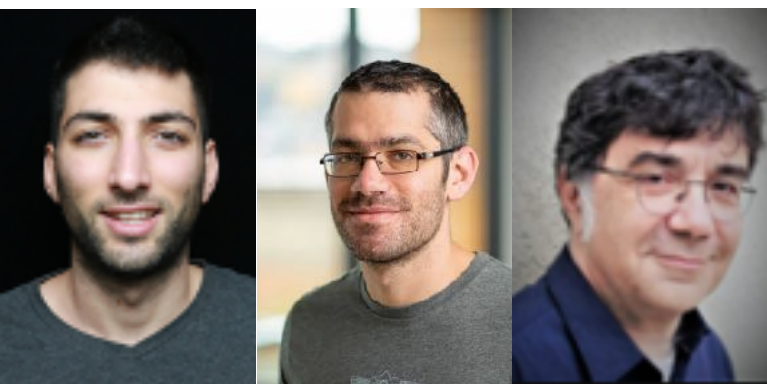
	WINOGr. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7_{0.4}	87.6_{0.6}

Data Efficient Masked Language Modeling for Vision and Language

Bitton, Stanovsky, Elhadad & S., Findings of EMNLP 2021

Baseline MLM

A tiger [MASK] eating the carrot



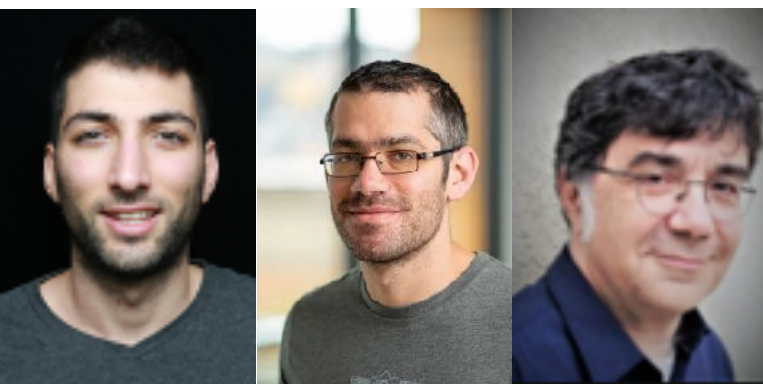
Data Efficient Masked Language Modeling for Vision and Language

Bitton, Stanovsky, Elhadad & S., Findings of EMNLP 2021

- Current practice: **randomly** mask some of the words in the sentence
 - Many of them are **stop words** and **punctuation**

Baseline MLM

A tiger [MASK] eating the carrot



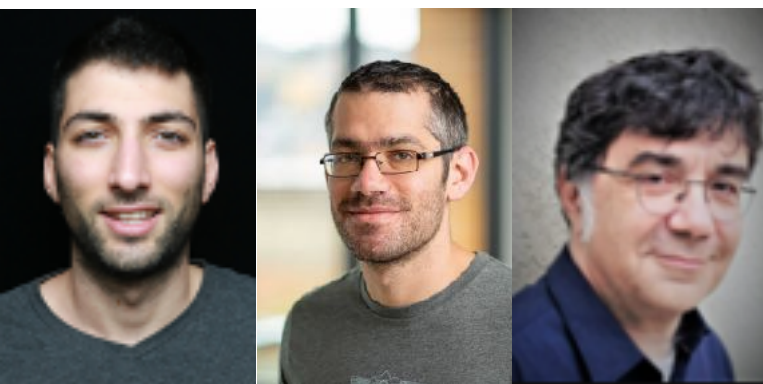
Data Efficient Masked Language Modeling for Vision and Language

Bitton, Stanovsky, Elhadad & S., Findings of EMNLP 2021

- Current practice: **randomly** mask some of the words in the sentence
 - Many of them are **stop words** and **punctuation**
- Our proposal: only mask **content words**

Baseline MLM

A tiger [MASK] eating the carrot



Data Efficient Masked Language Modeling for Vision and Language

Bitton, Stanovsky, Elhadad & S., Findings of EMNLP 2021

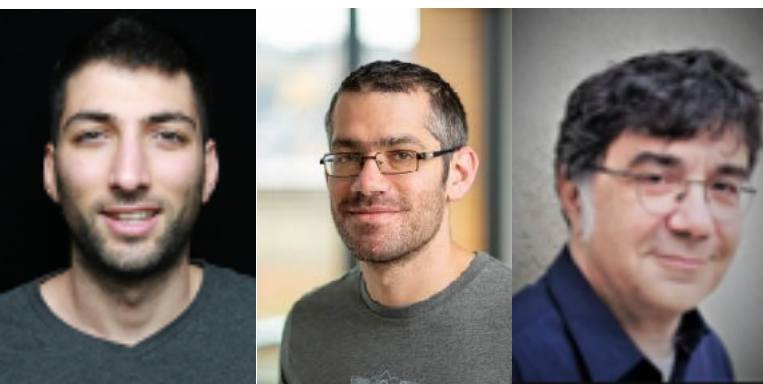
- Current practice: **randomly** mask some of the words in the sentence
 - Many of them are **stop words** and **punctuation**
- Our proposal: only mask **content words**

Baseline MLM

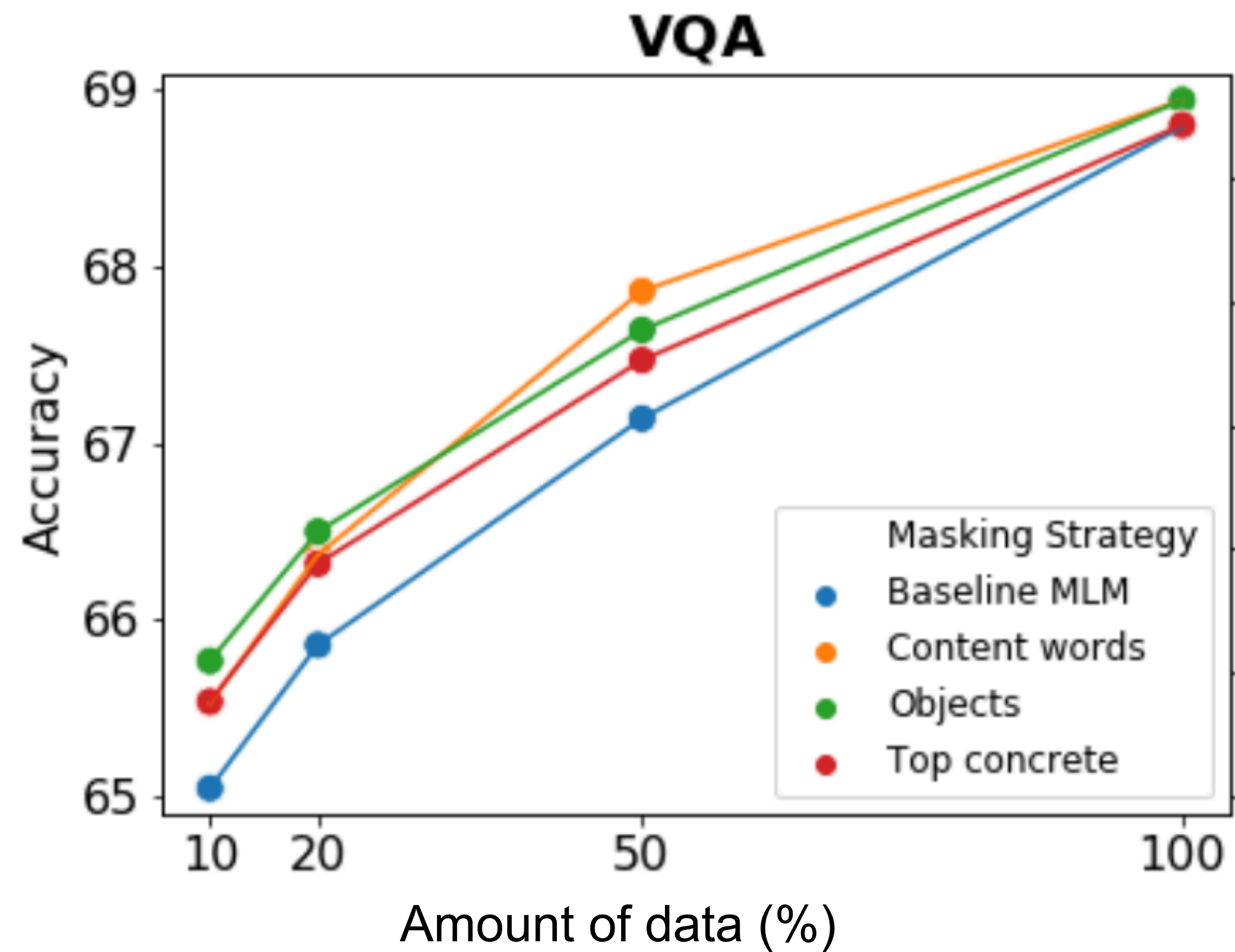
A tiger [MASK] eating the carrot

Our method

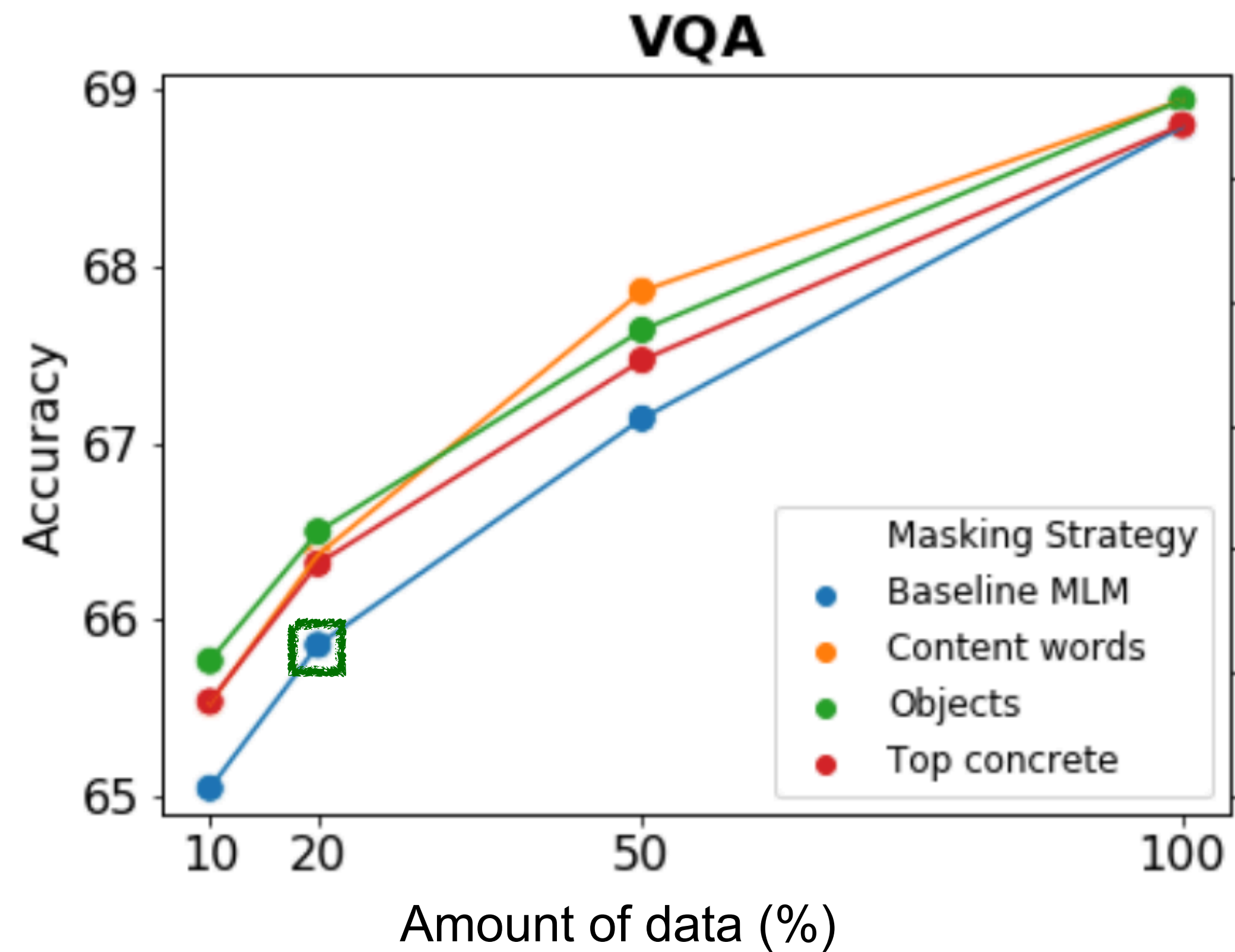
A [MASK] is eating the carrot



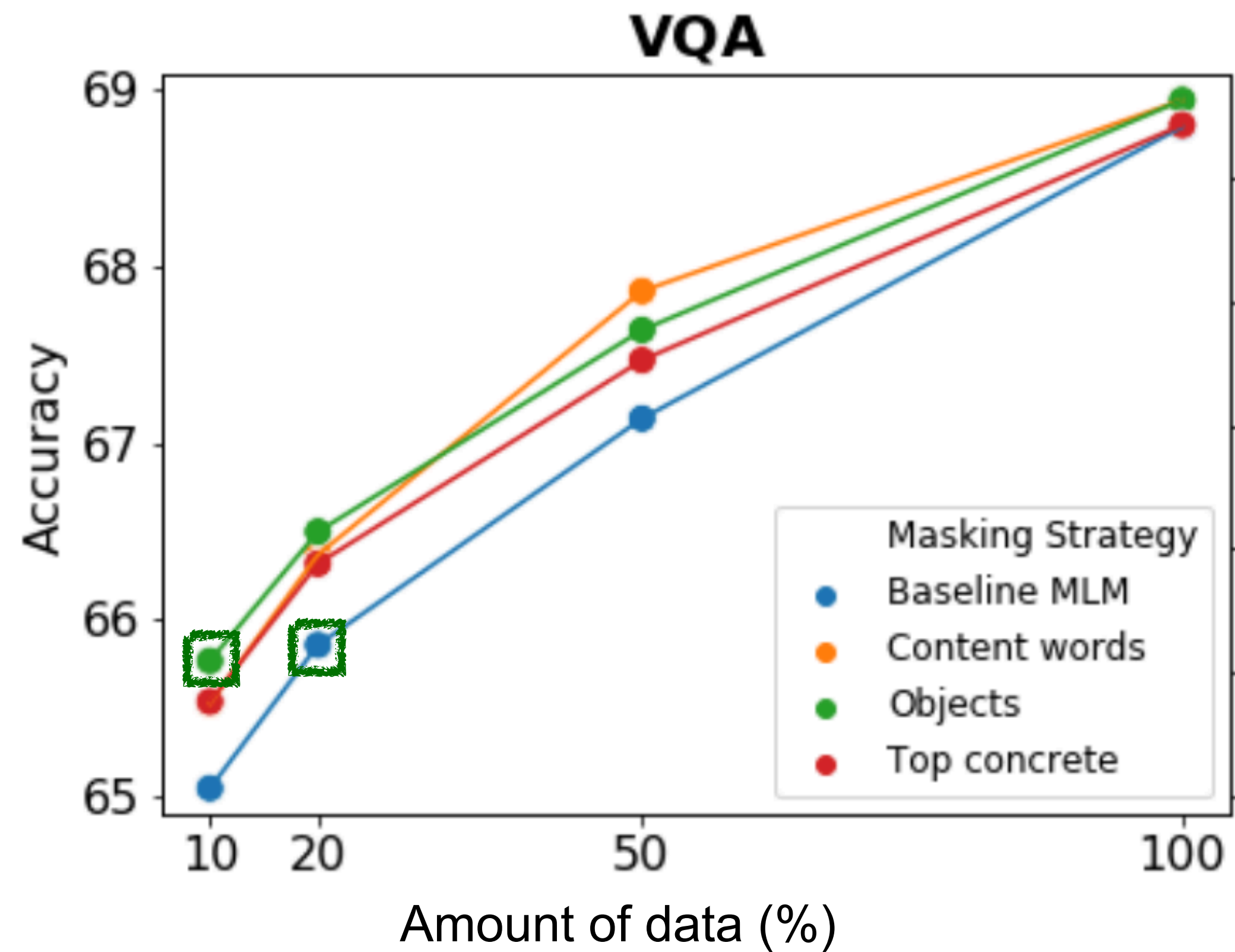
Data Efficient Training



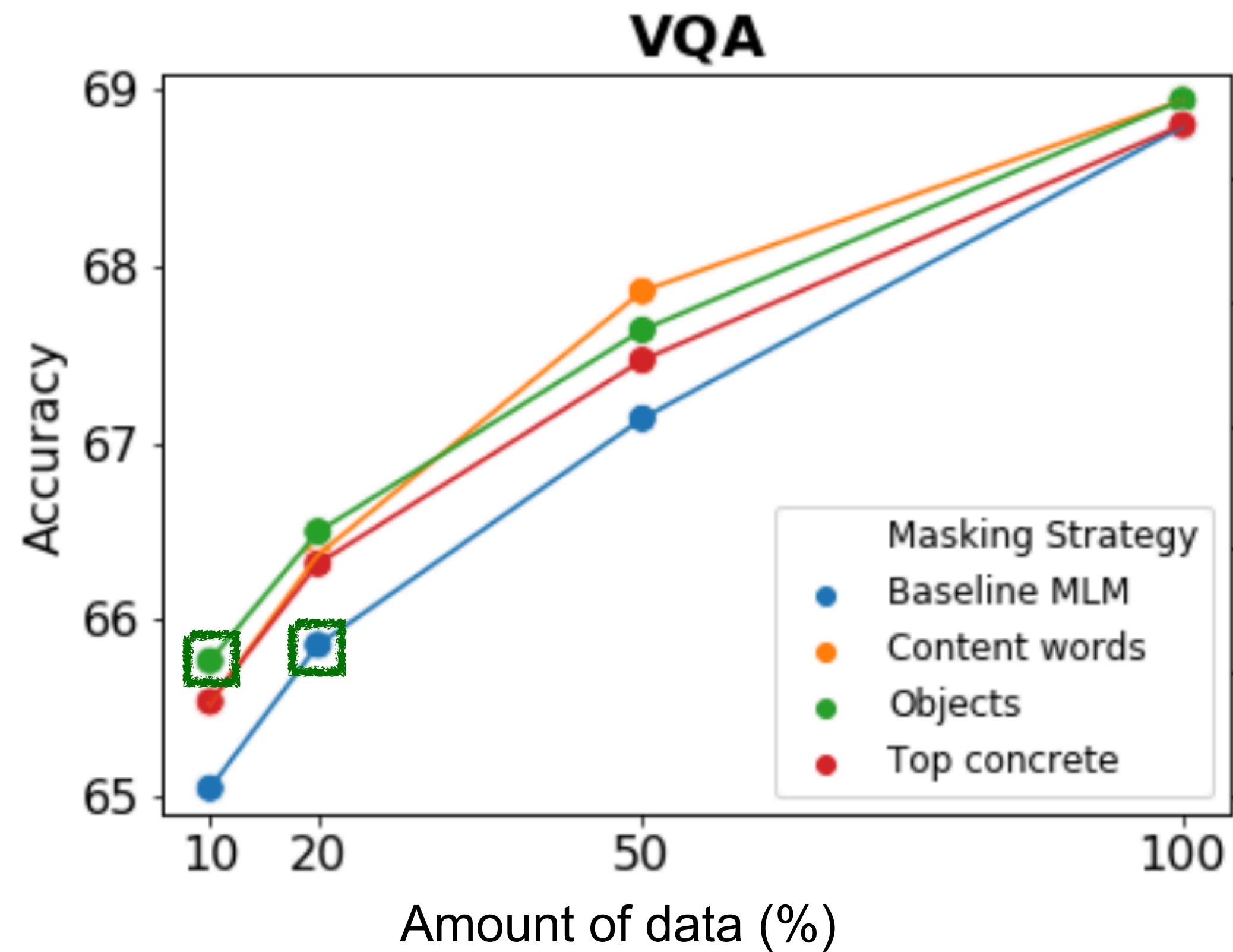
Data Efficient Training



Data Efficient Training



Data Efficient Training



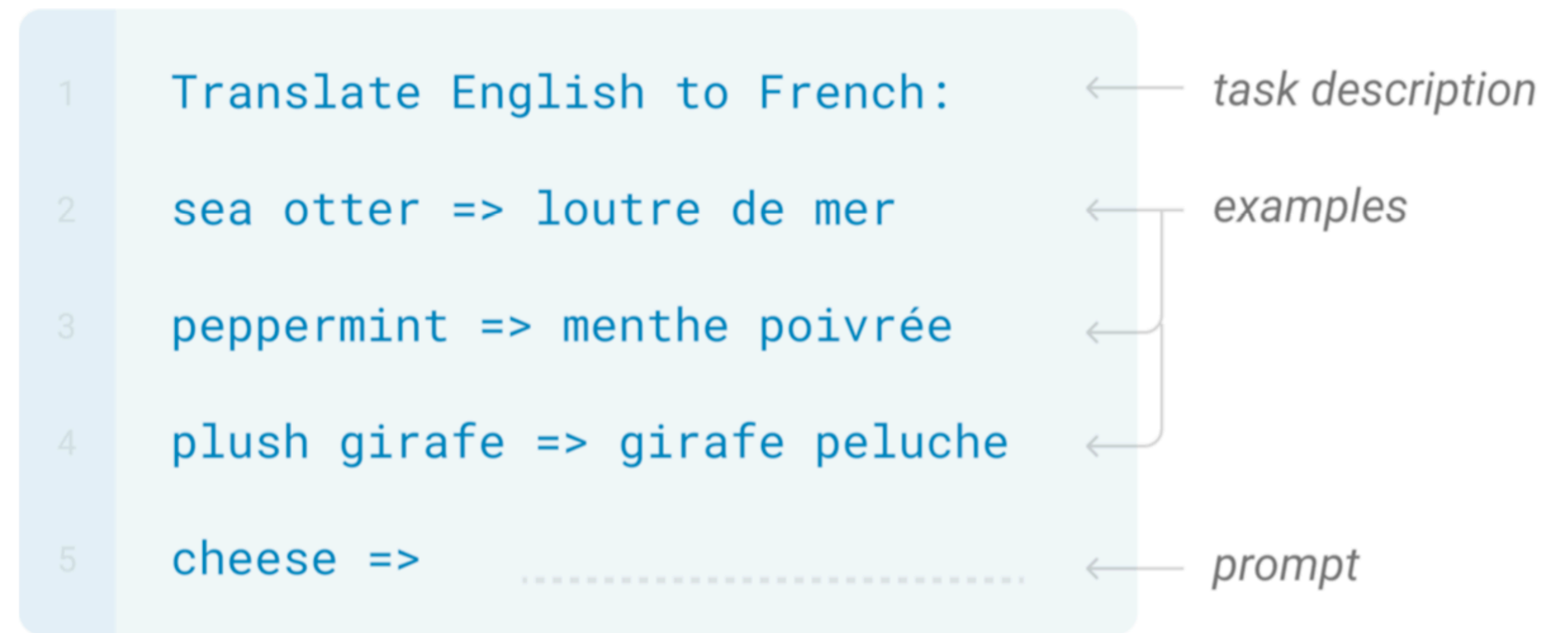
Similar accuracy, twice as fast

Few-shot Learning

- Only use a handful of examples to train a model

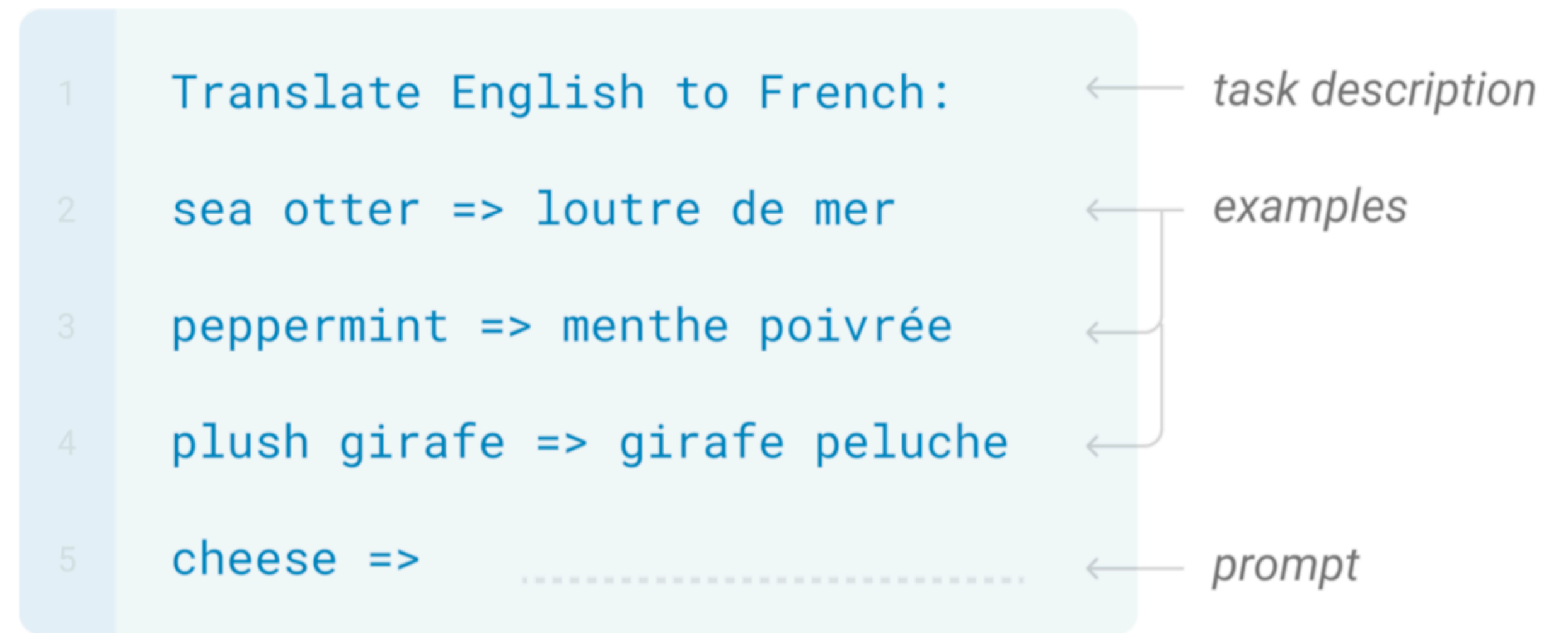
Few-shot Learning

- Only use a handful of examples to train a model
- Prompting
 - Brown et al. (2020), Schick & Schütze, 2021)



Few-shot Learning

- Only use a handful of examples to train a model
- Prompting
 - Brown et al. (2020), Schick & Schütze, 2021)
- Non-prompting methods
 - Mahabadi et al. (2022)



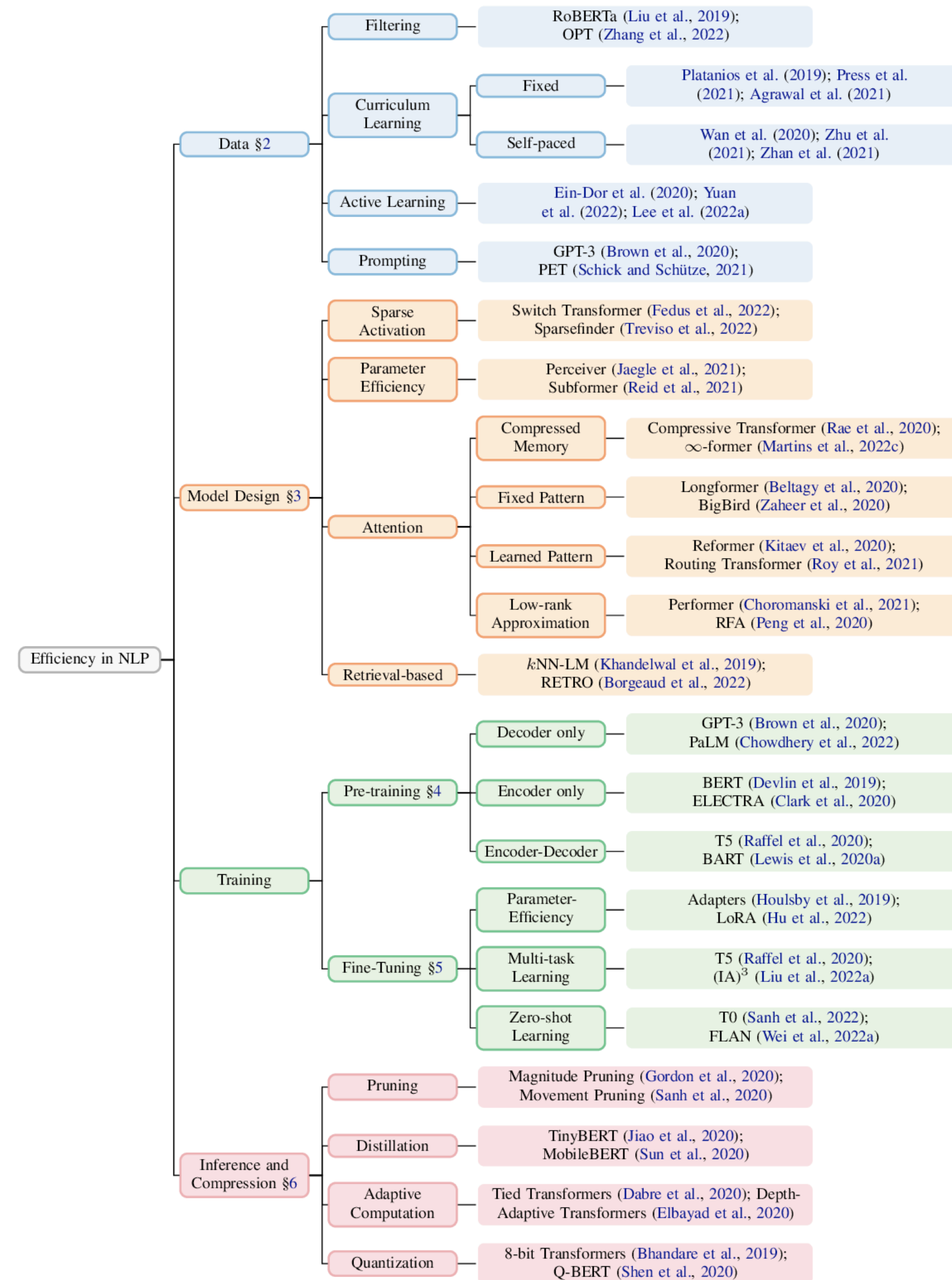
Data Efficiency

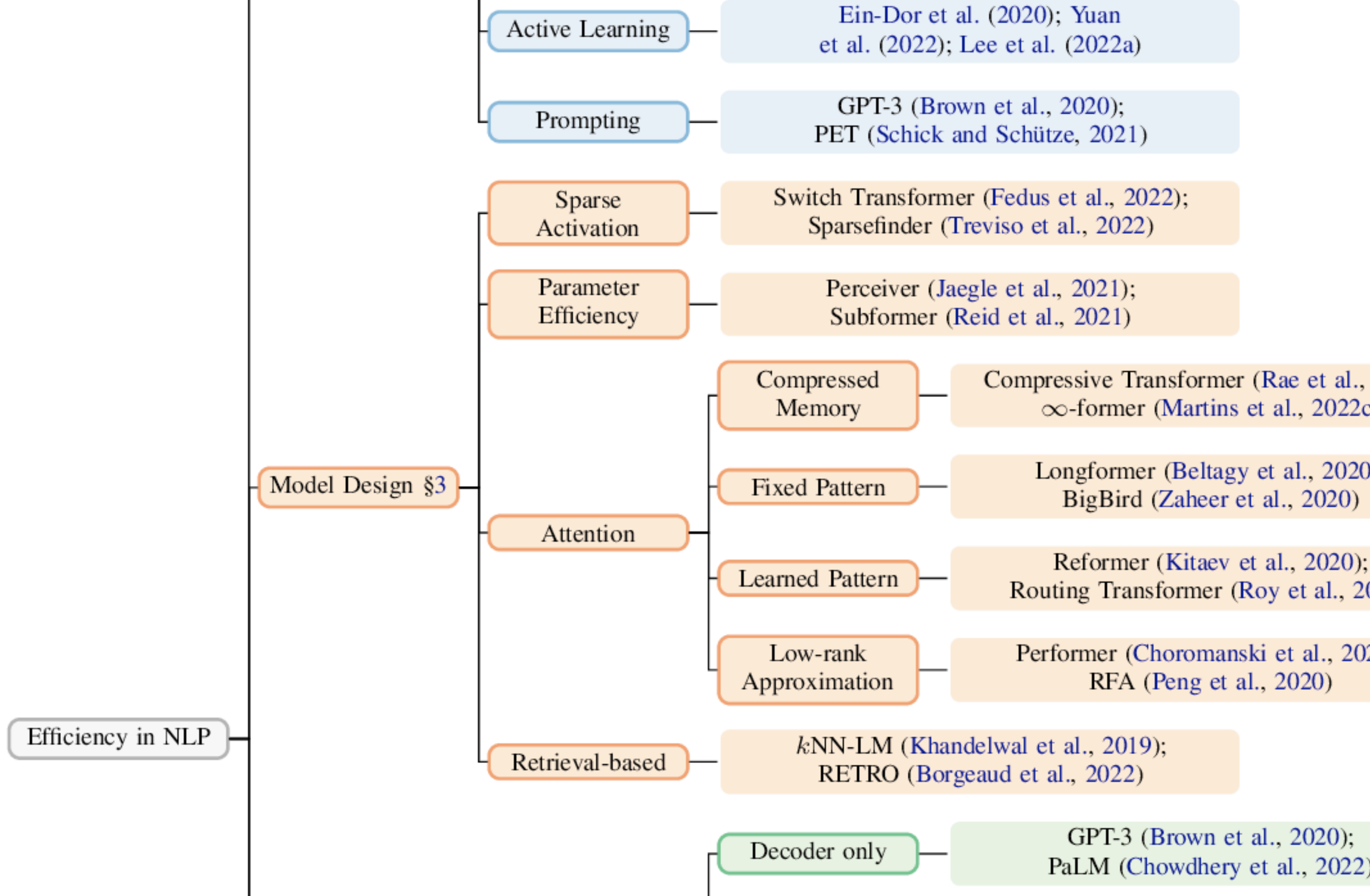
Open Questions



- Do we really need massive web-scale data to train our models?
 - Can we get along with less?
 - Sorscher et al. (2022)

Efficient Methods in NLP



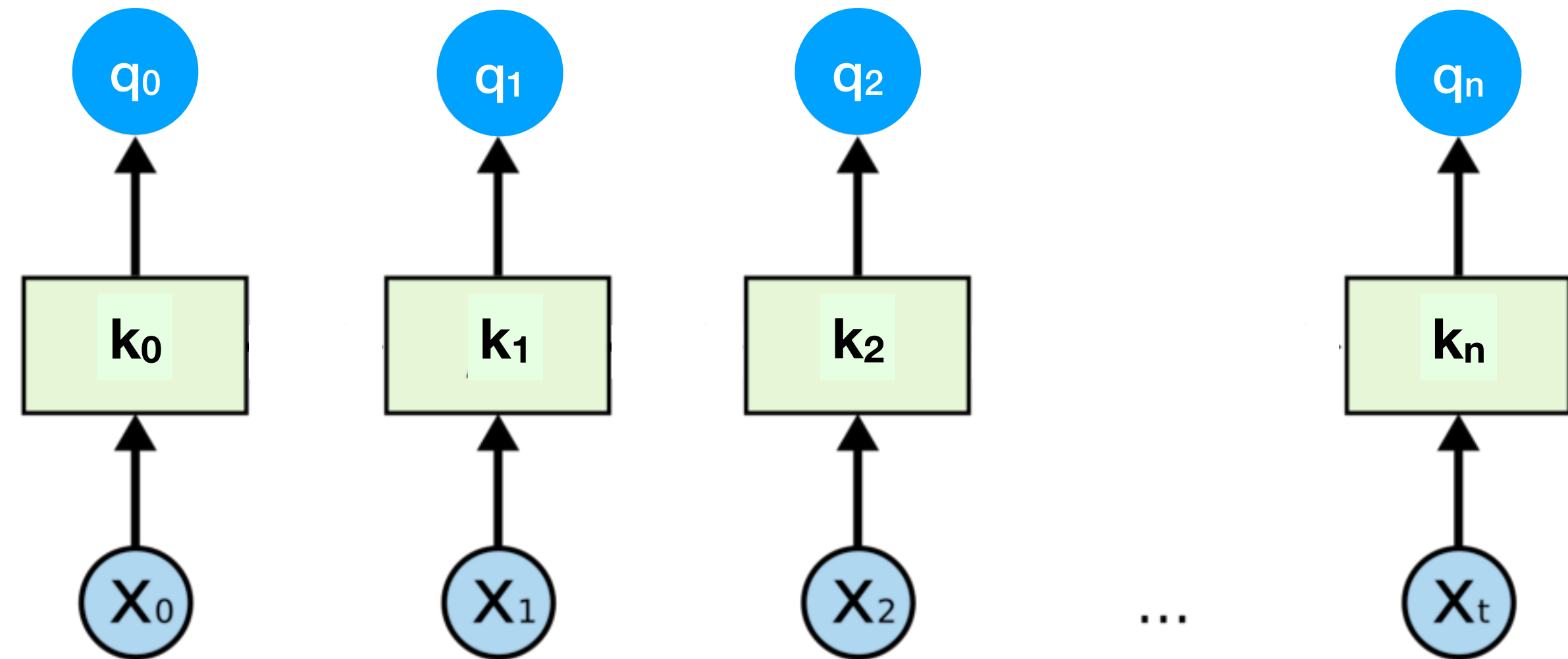


Transformers

Vaswani et al., 2017



- **The** method for text representation
 - Also for vision, speech, combio, ...

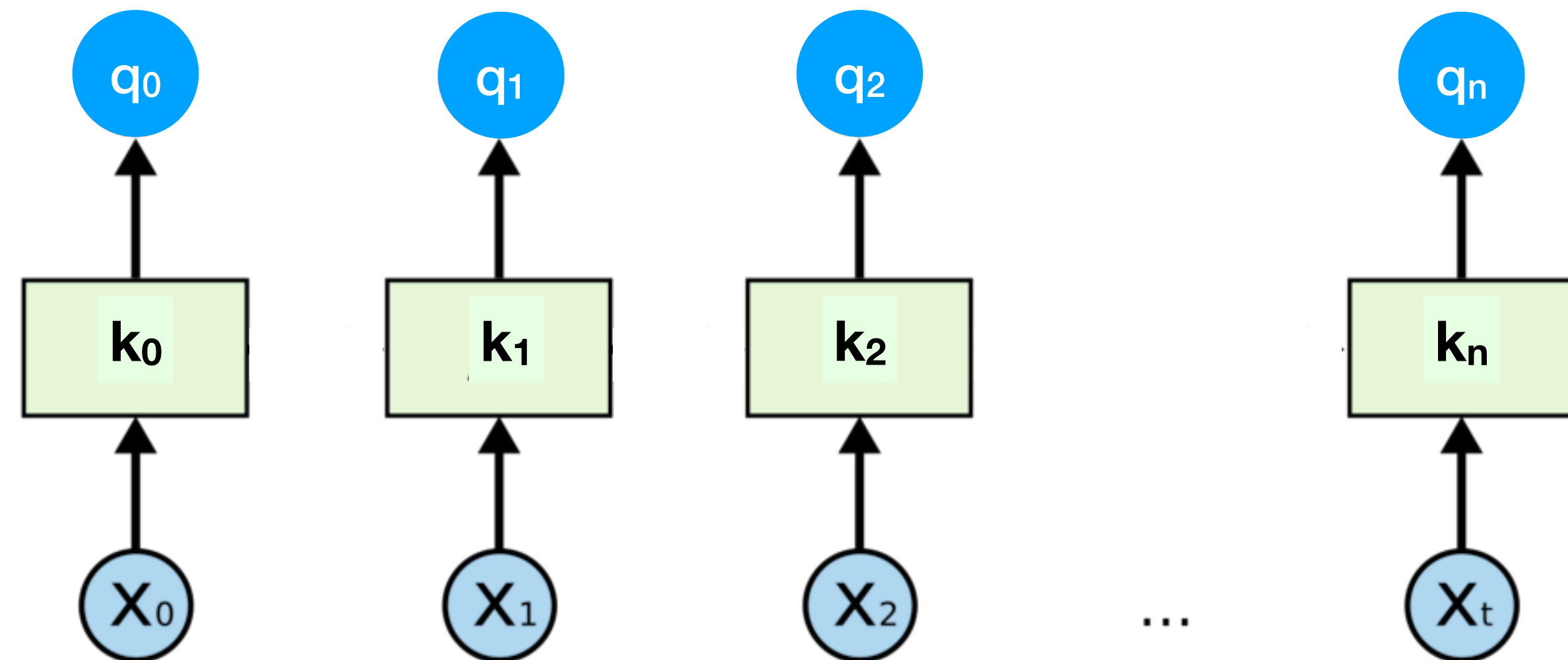


Transformers

Vaswani et al., 2017



- **The** method for text representation
 - Also for vision, speech, combio, ...
- Each word attends to all other words

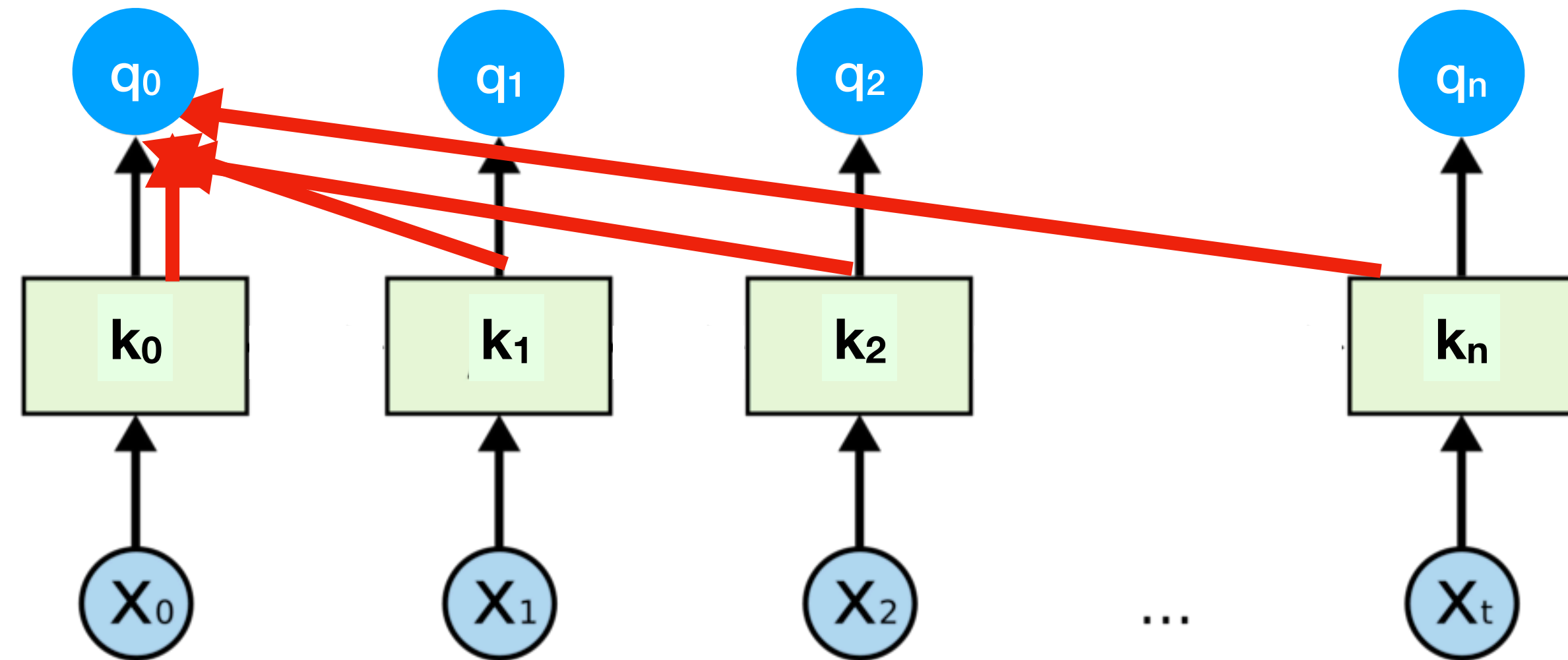


Transformers

Vaswani et al., 2017



- **The** method for text representation
 - Also for vision, speech, combio, ...
- Each word attends to all other words

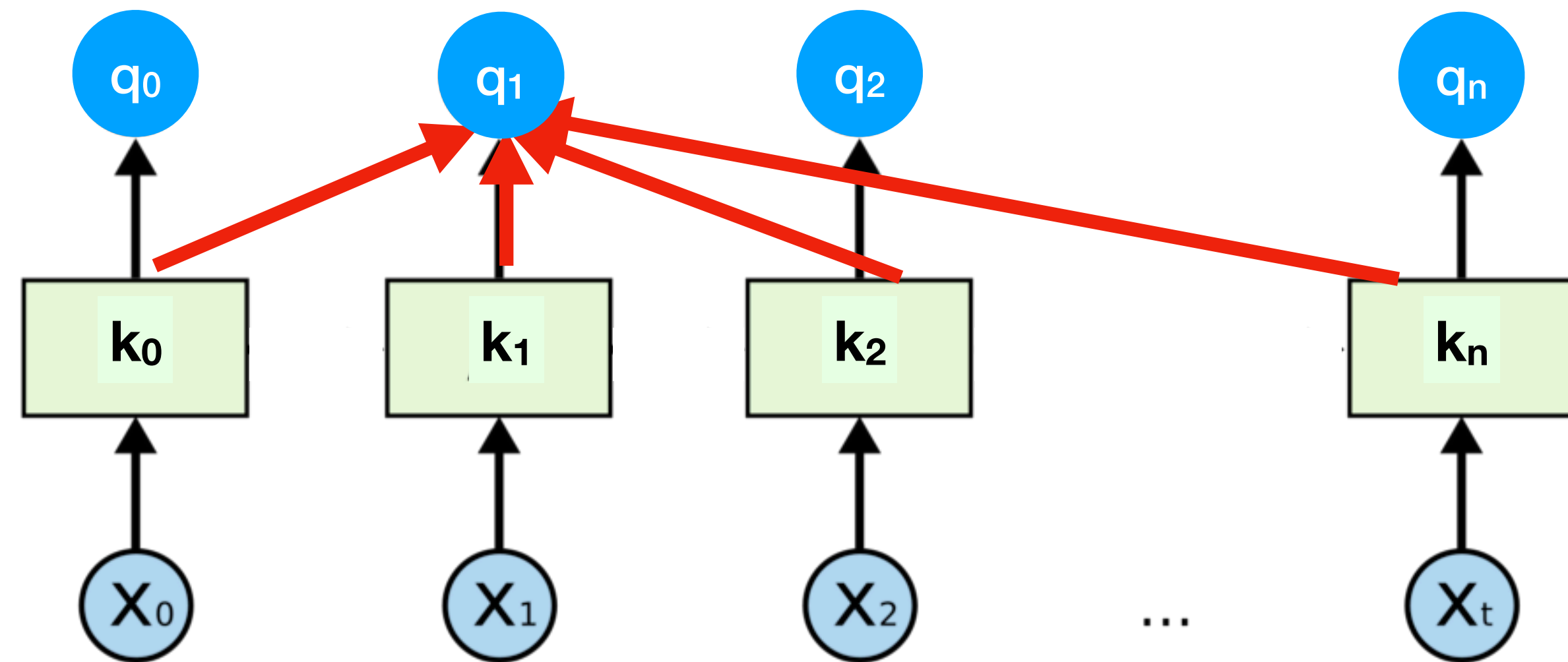


Transformers

Vaswani et al., 2017



- **The** method for text representation
 - Also for vision, speech, combio, ...
- Each word attends to all other words

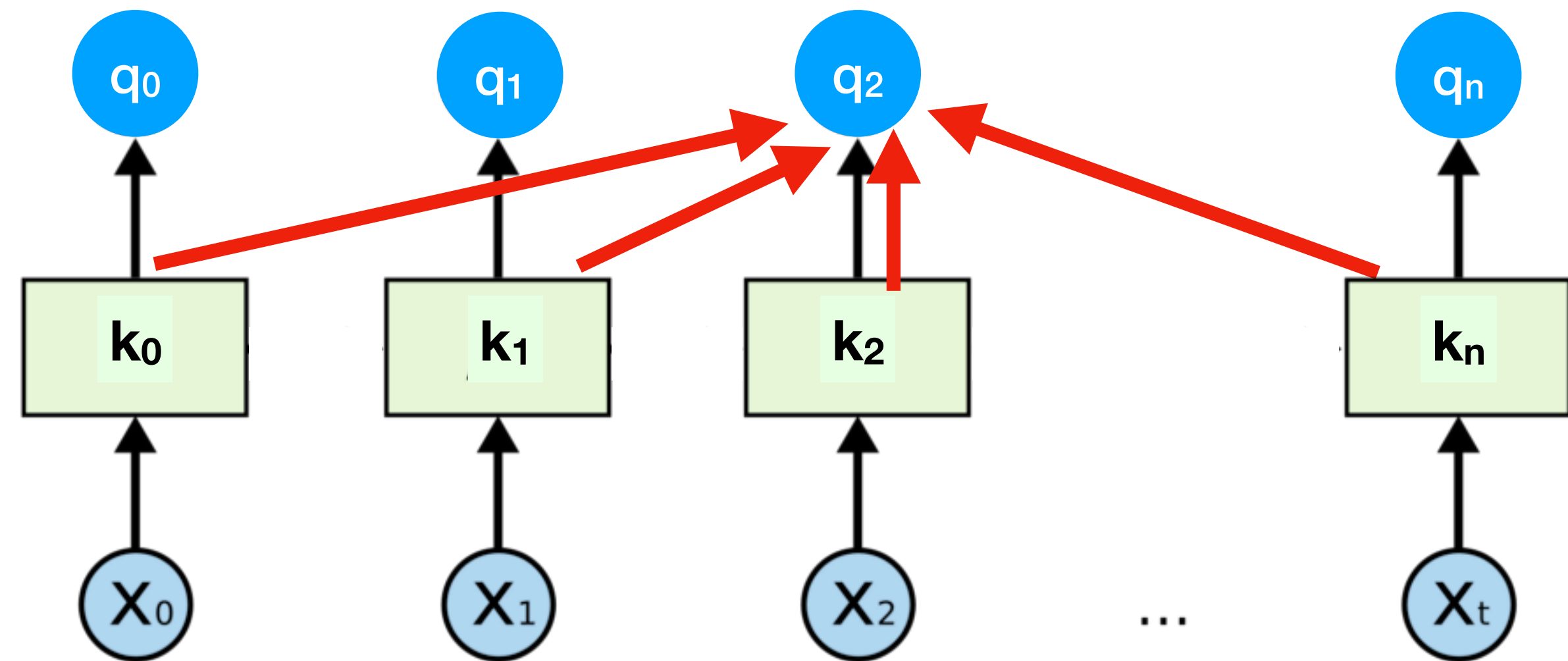


Transformers

Vaswani et al., 2017



- **The** method for text representation
 - Also for vision, speech, combio, ...
- Each word attends to all other words

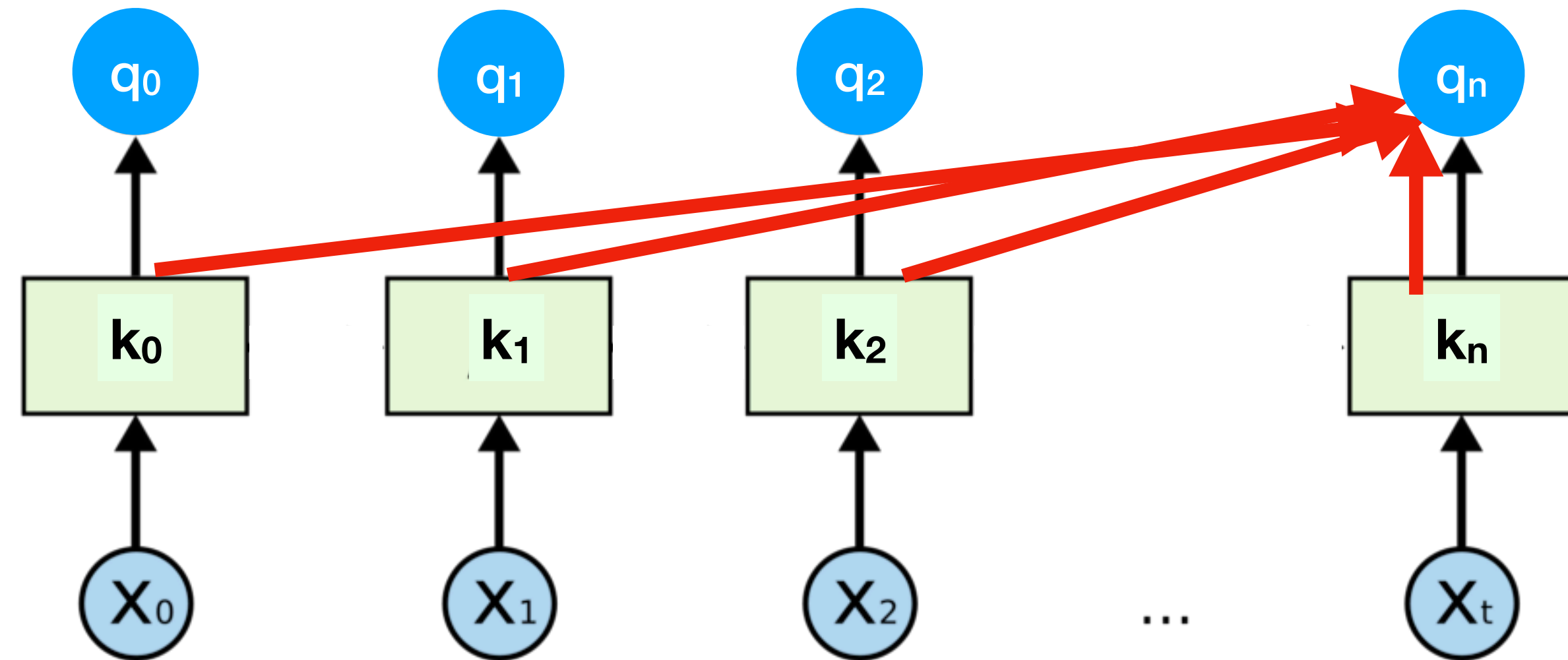


Transformers

Vaswani et al., 2017



- **The** method for text representation
 - Also for vision, speech, combio, ...
- Each word attends to all other words

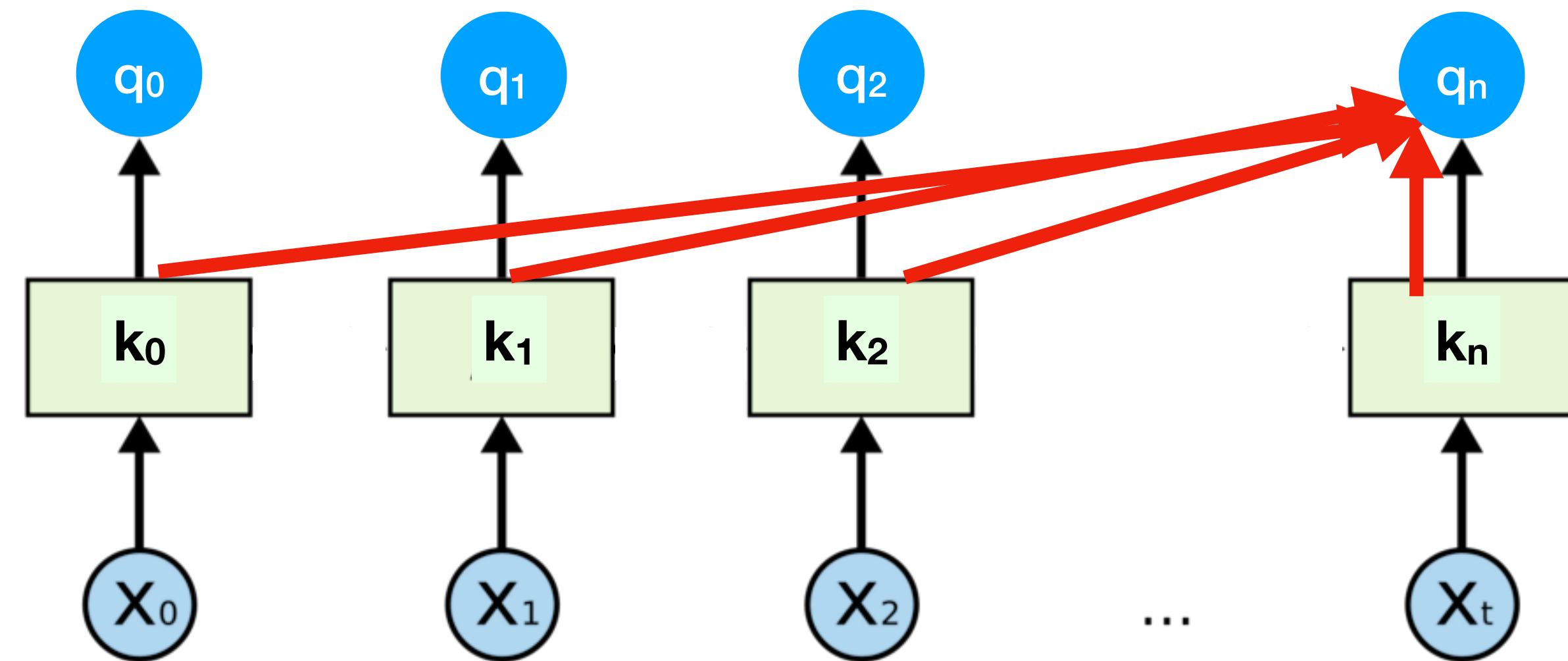


Transformers

Vaswani et al., 2017



- **The** method for text representation
 - Also for vision, speech, combio, ...
- Each word attends to all other words
- $O(n^2)$ complexity in the sentence length n

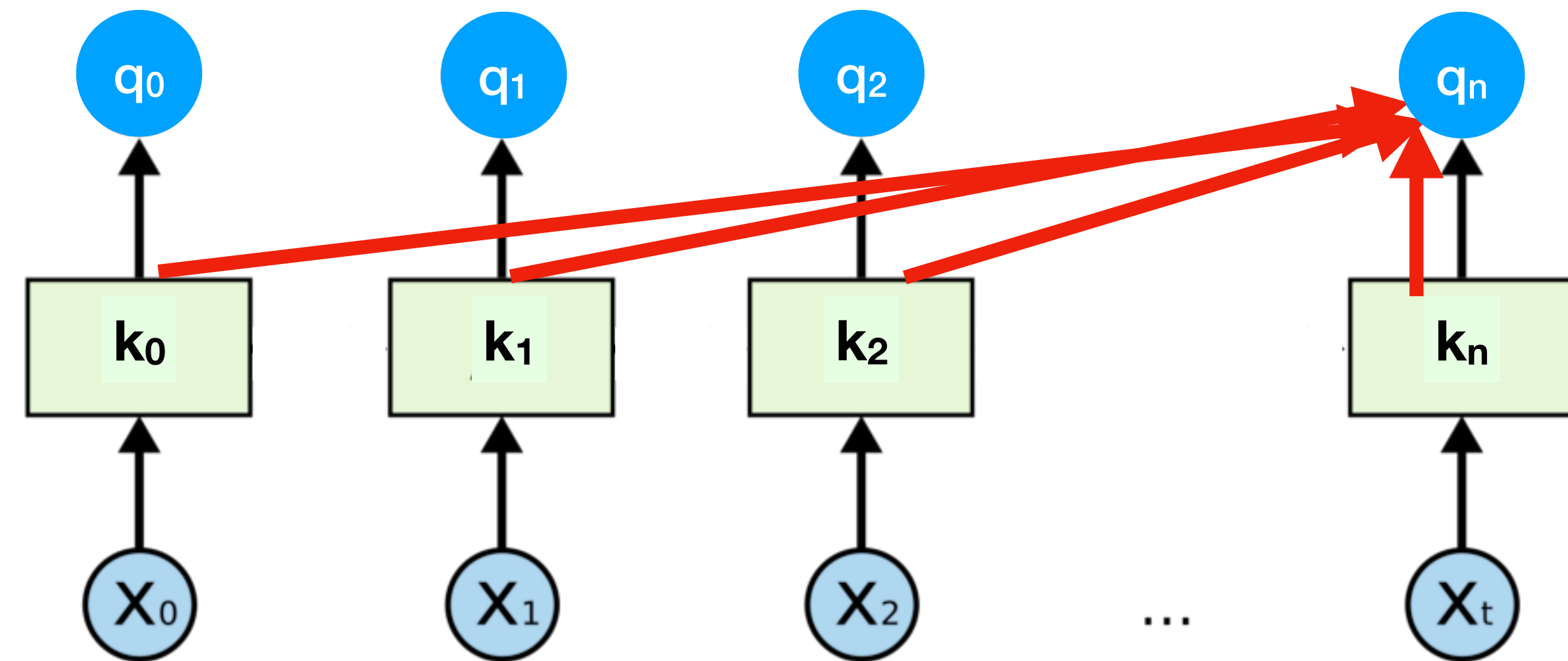


Transformers

Vaswani et al., 2017



- **The** method for text representation
 - Also for vision, speech, combio, ...
- Each word attends to all other words
- $O(n^2)$ complexity in the sentence length n
- Fatal for long sequences
 - Books, articles, etc.



Random Feature Attention

Peng, Pappas, Yogatama, S., Smith, & Kong, ICLR 2021

spotlight presentation

- **Key idea:** approximate the attention function using random Fourier features
 - Rahimi and Recht (2007)

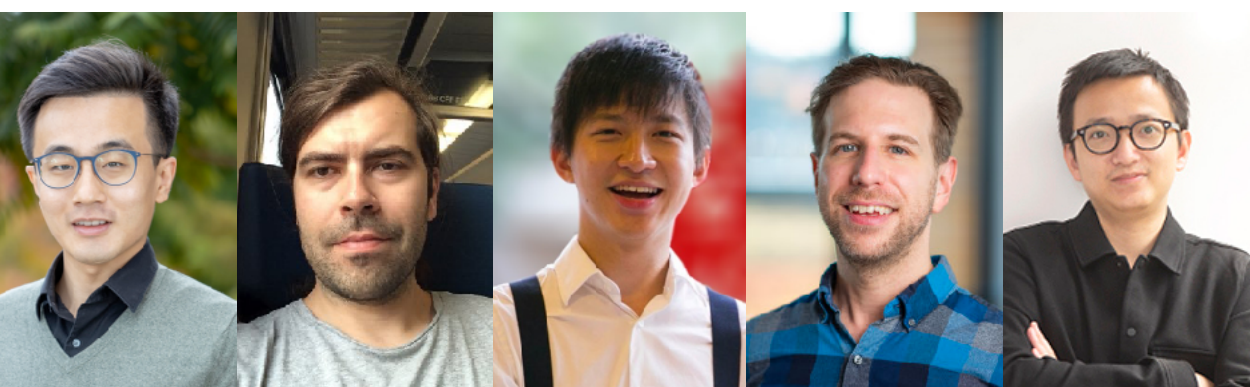


Random Feature Attention

Peng, Pappas, Yogatama, **S.**, Smith, & Kong, ICLR 2021

spotlight presentation

- **Key idea:** approximate the attention function using random Fourier features
 - Rahimi and Recht (2007)
- Some math

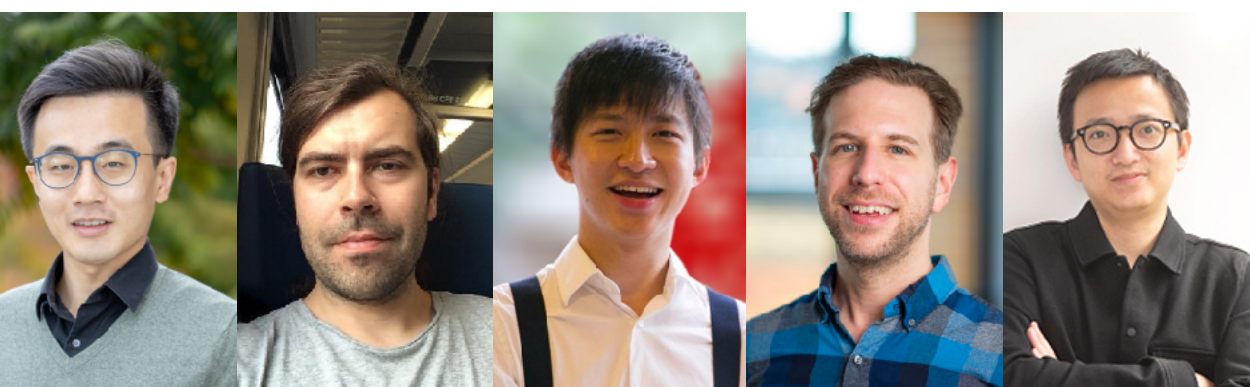


Random Feature Attention

Peng, Pappas, Yogatama, **S.**, Smith, & Kong, ICLR 2021

spotlight presentation

- **Key idea:** approximate the attention function using random Fourier features
 - Rahimi and Recht (2007)
- Some math

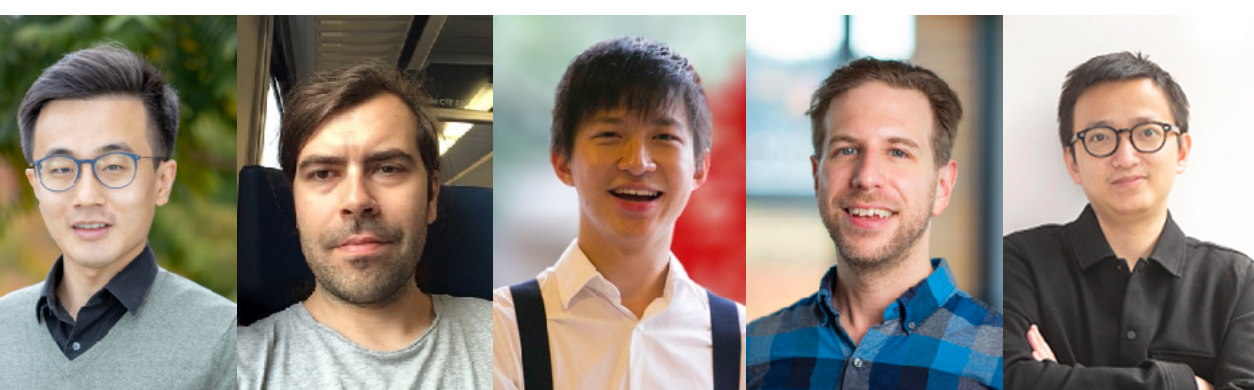


Random Feature Attention

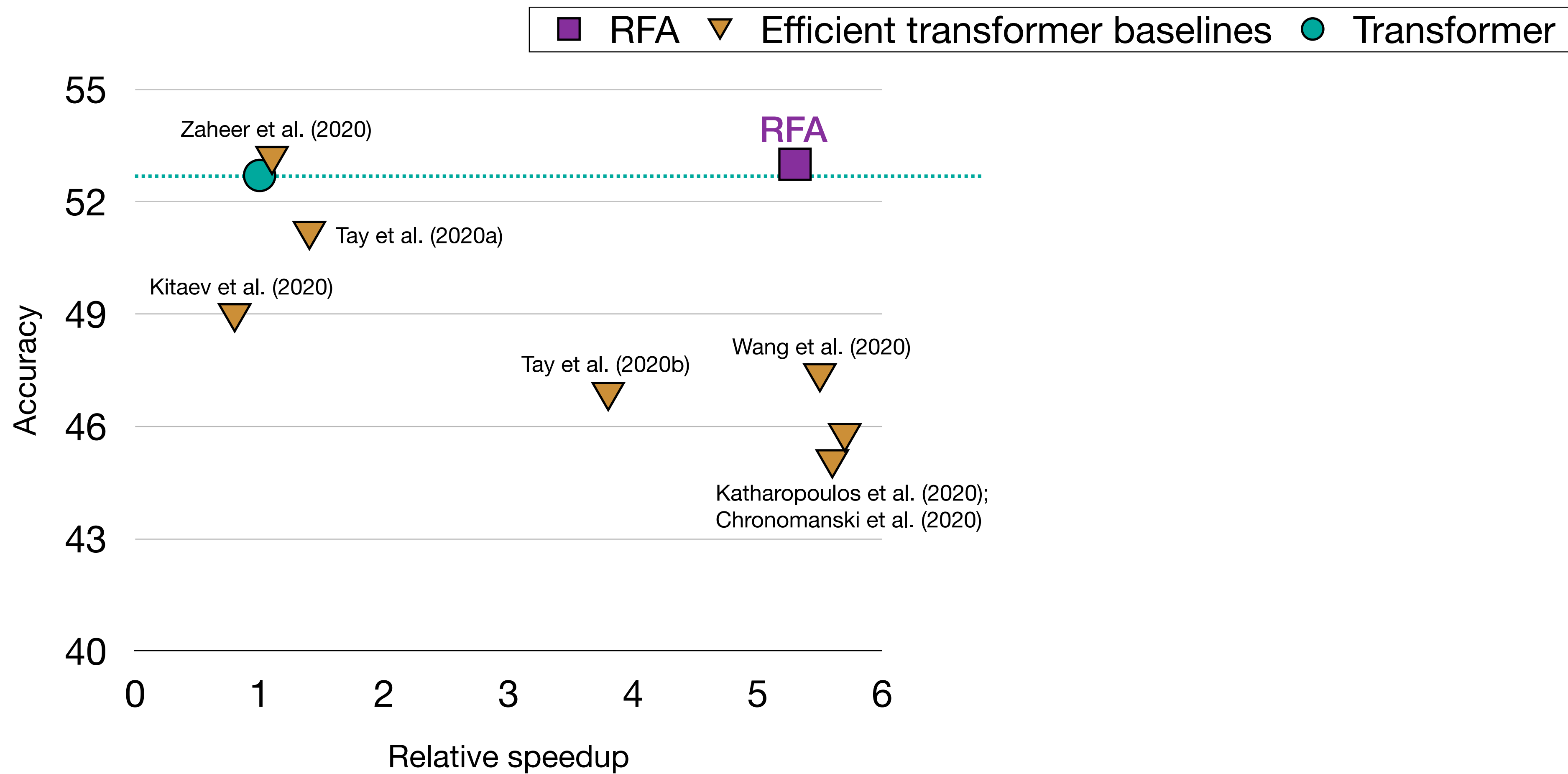
Peng, Pappas, Yogatama, S., Smith, & Kong, ICLR 2021

spotlight presentation

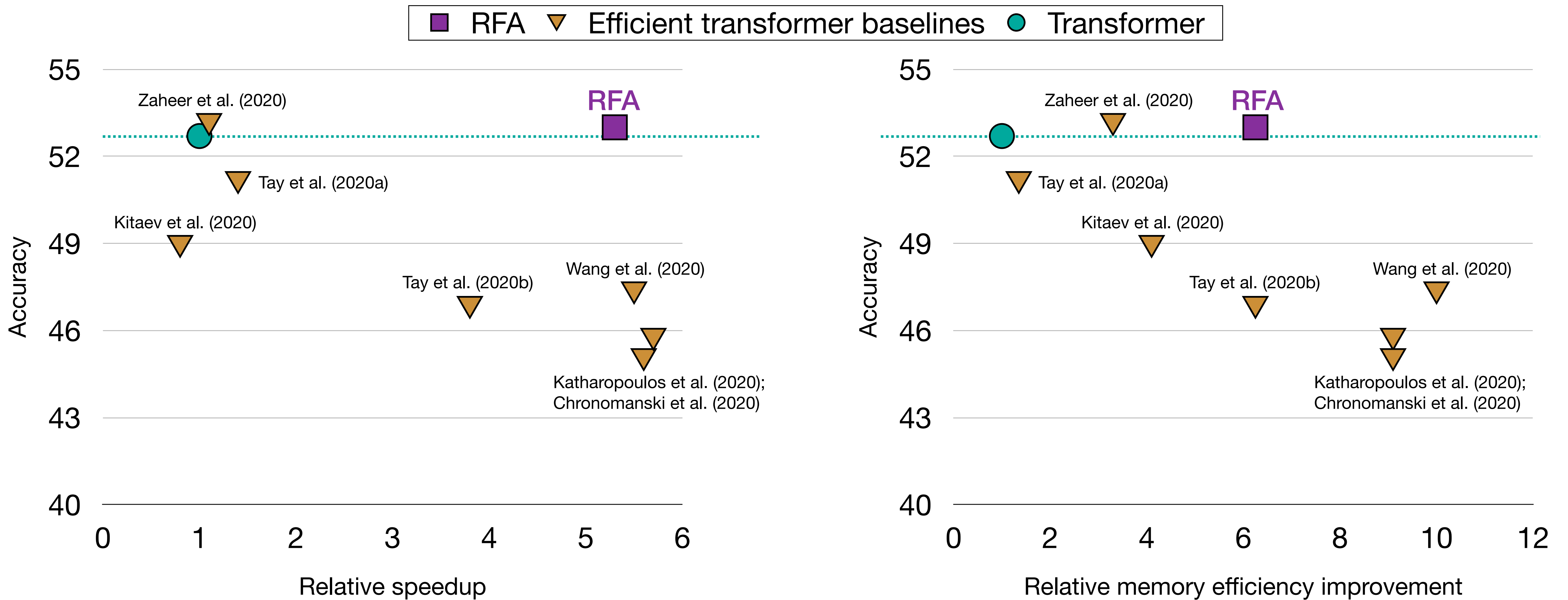
- **Key idea:** approximate the attention function using random Fourier features
 - Rahimi and Recht (2007)
- Some math
- Linear runtime and memory requirements



Better Efficiency-Accuracy Tradeoff



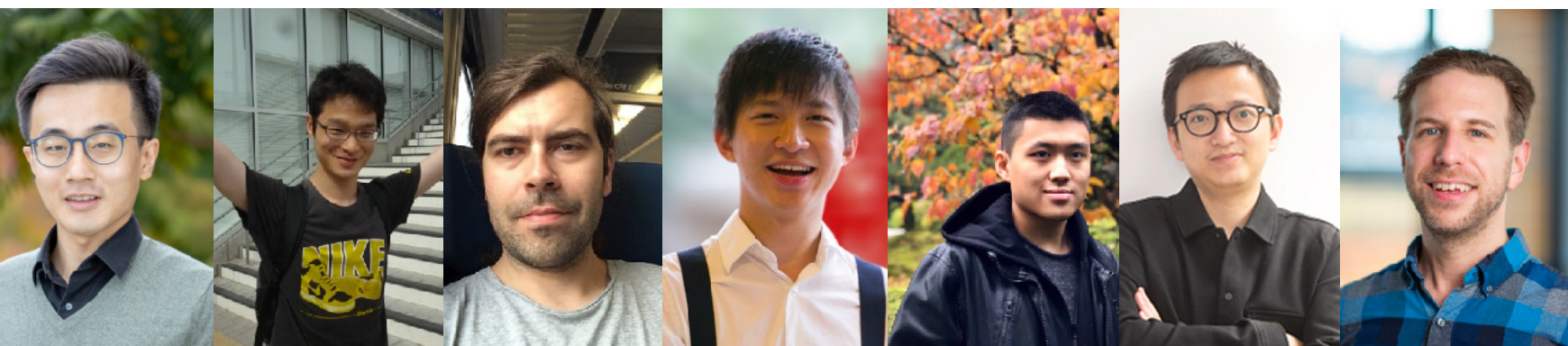
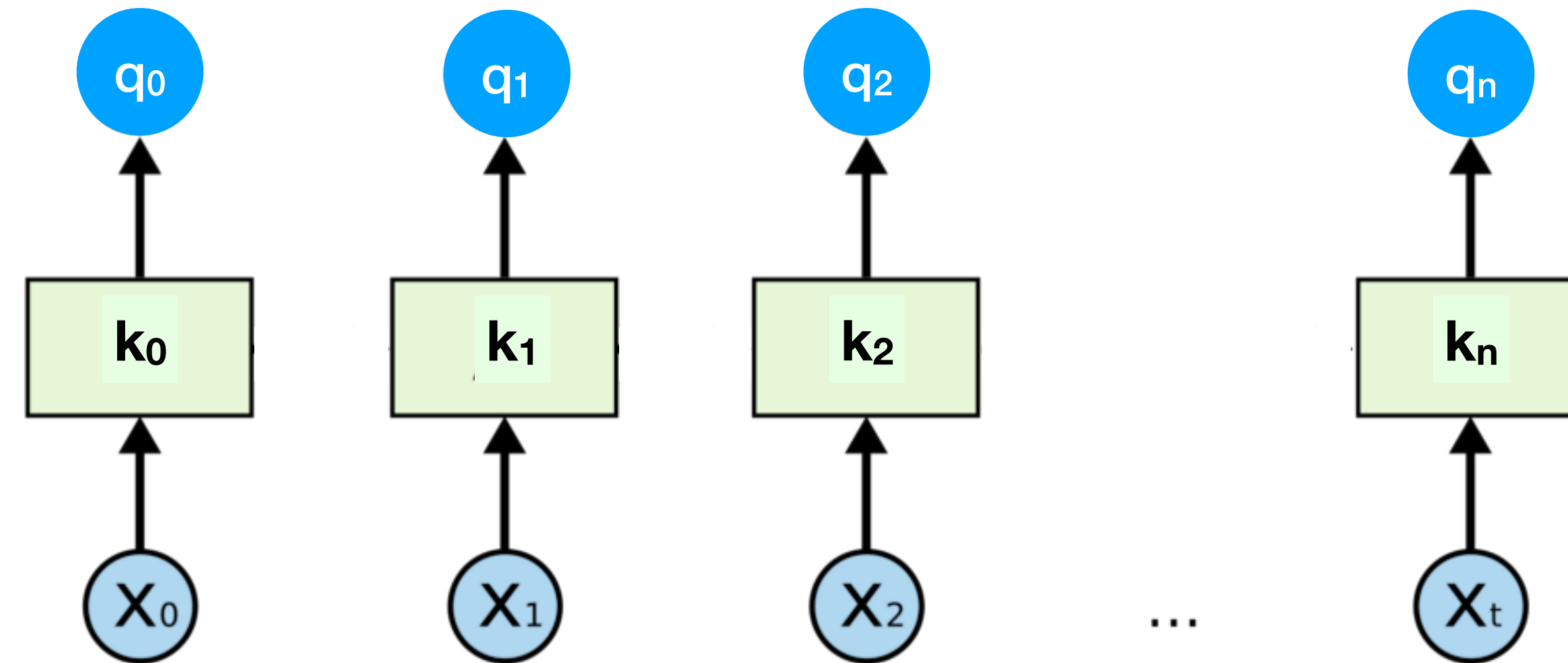
Better Efficiency-Accuracy Tradeoff



ABC: Attention with Bounded-memory Control

Peng, Kasai, Pappas, Yogatama, Wu, Kong, S. & Smith, ACL 2022

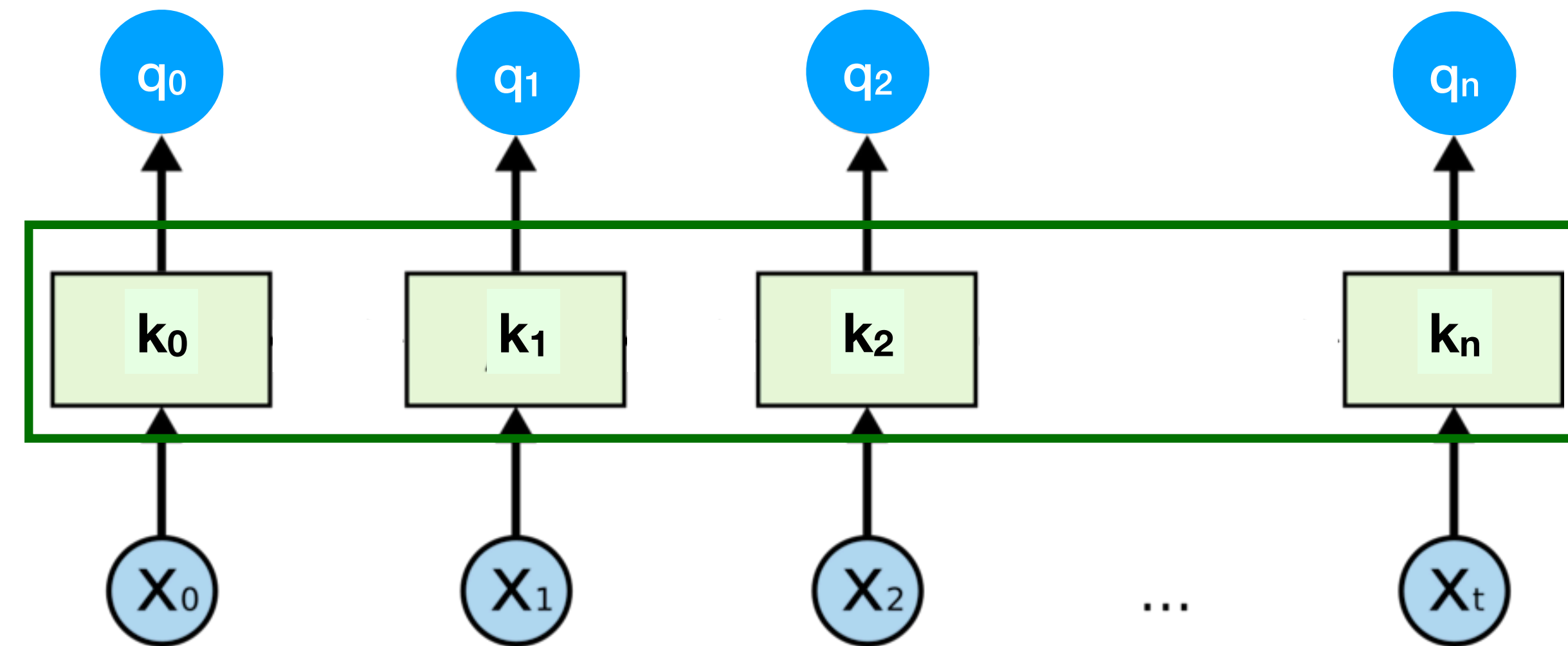
- **Key intuition:** treat the sentence as memory of size n



ABC: Attention with Bounded-memory Control

Peng, Kasai, Pappas, Yogatama, Wu, Kong, S. & Smith, ACL 2022

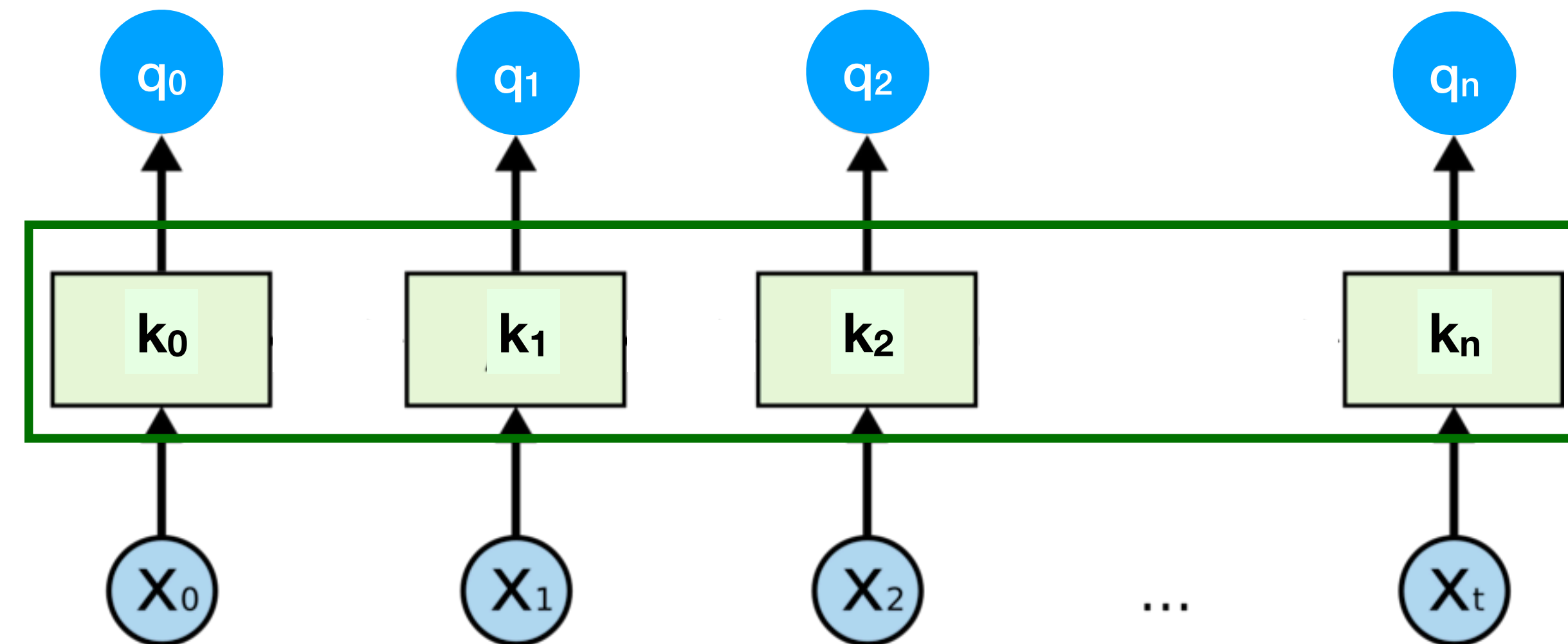
- **Key intuition:** treat the sentence as memory of size n



ABC: Attention with Bounded-memory Control

Peng, Kasai, Pappas, Yogatama, Wu, Kong, S. & Smith, ACL 2022

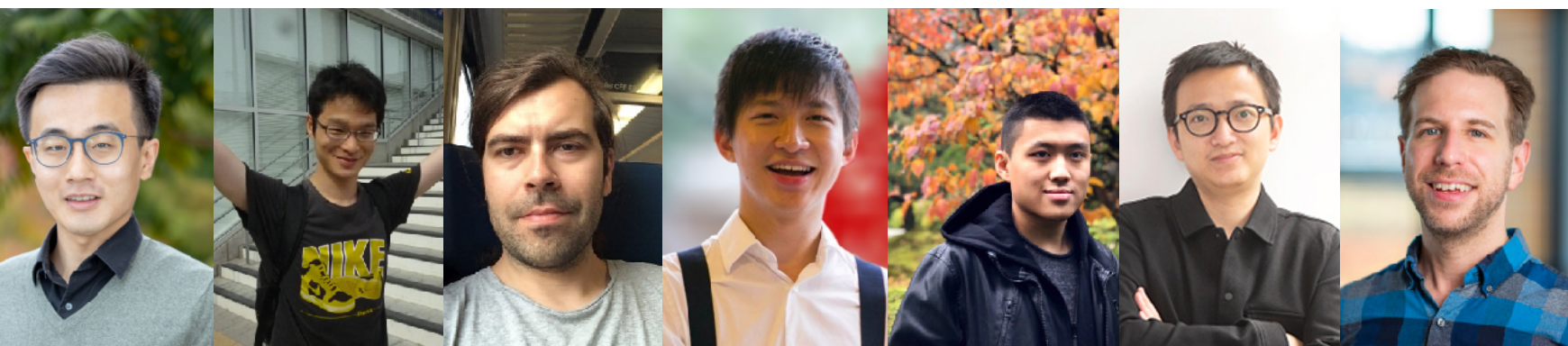
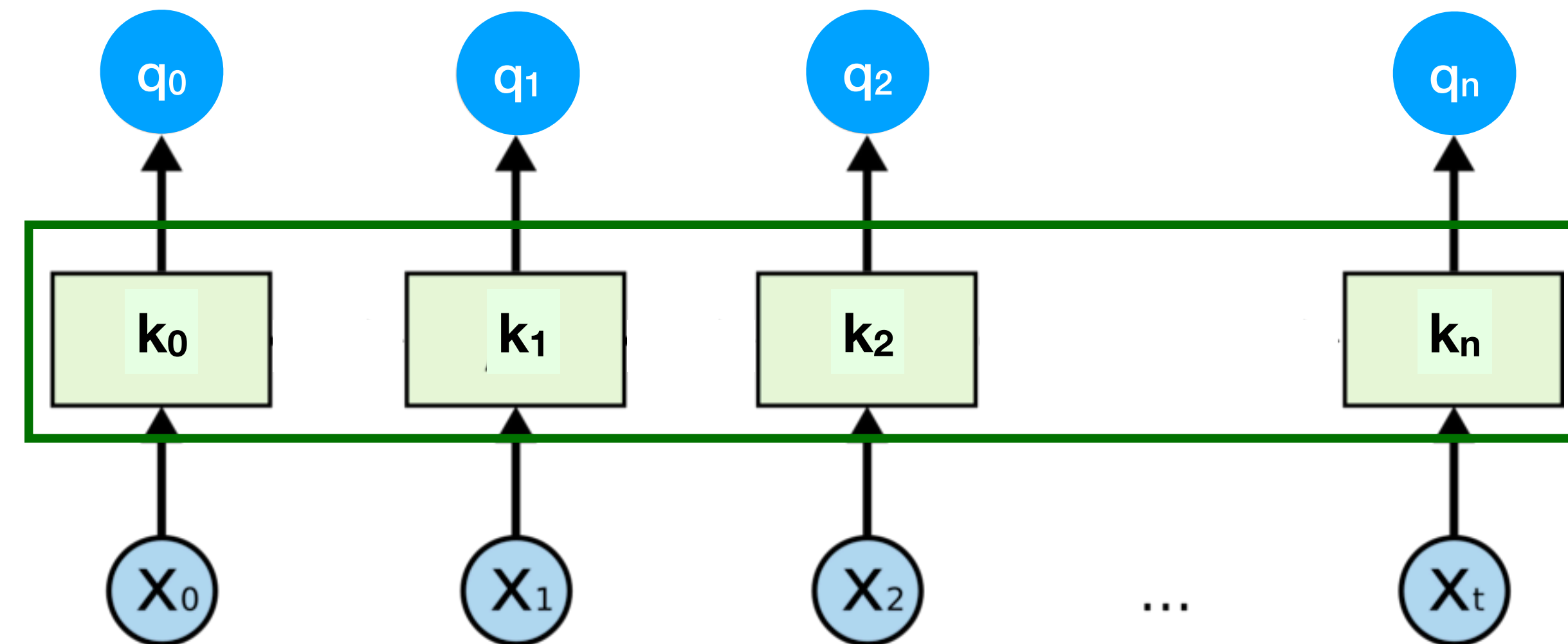
- **Key intuition:** treat the sentence as memory of size n
- **Key idea:** replace this memory with a fixed size memory of (fixed) size $k \ll n$
 - Instead of attending n tokens, each word attends to k tokens



ABC: Attention with Bounded-memory Control

Peng, Kasai, Pappas, Yogatama, Wu, Kong, S. & Smith, ACL 2022

- **Key intuition:** treat the sentence as memory of size n
- **Key idea:** replace this memory with a fixed size memory of (fixed) size $k \ll n$
 - Instead of attending n tokens, each word attends to k tokens
- Overall complexity linear in n
 - With constant k



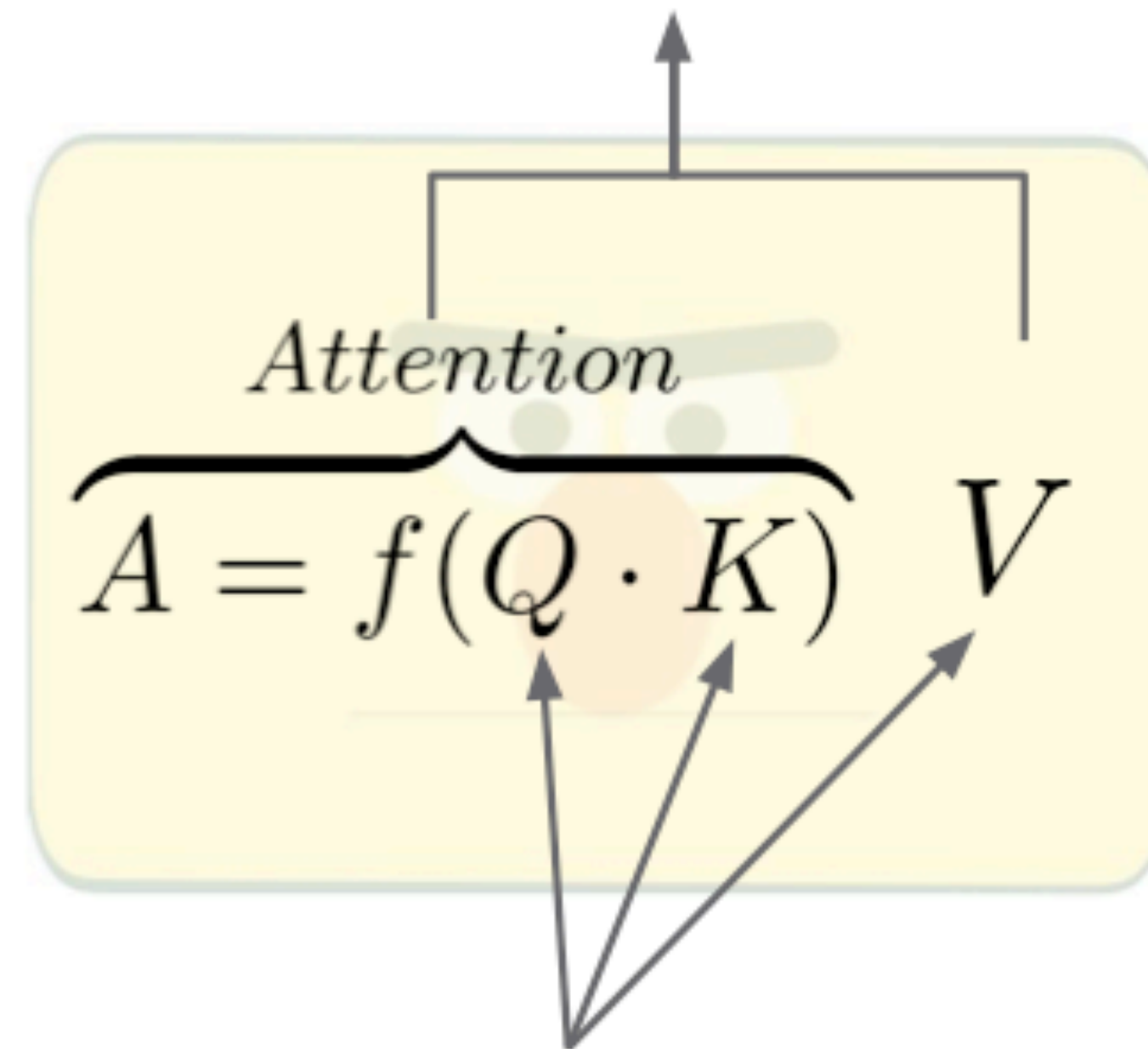
ABC Results

Speed

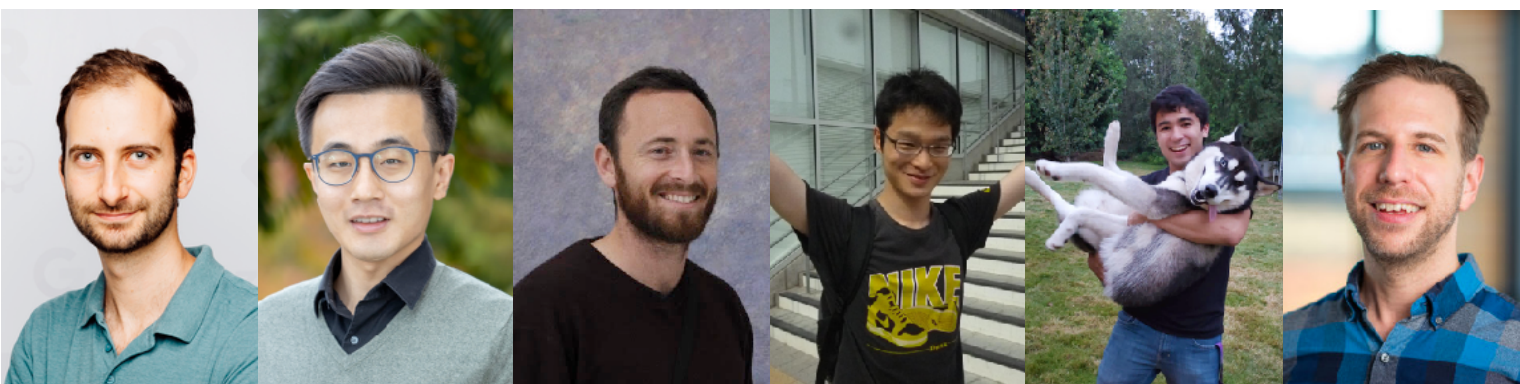
Memory

How Much Does Attention Actually Attend?

Hassid, Peng, Rotem, Kasai, Montero, Smith & S., Findings of EMNLP 2022

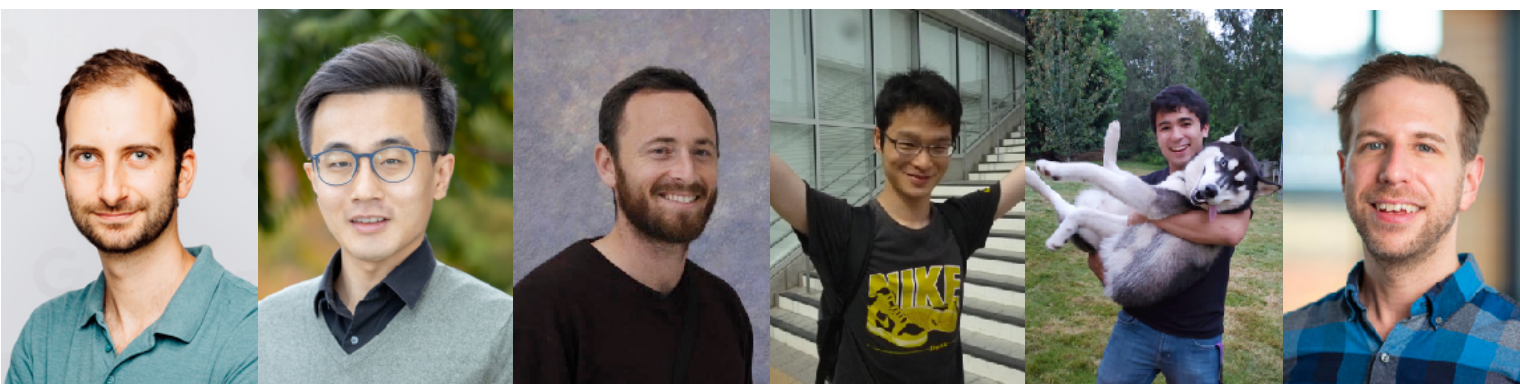
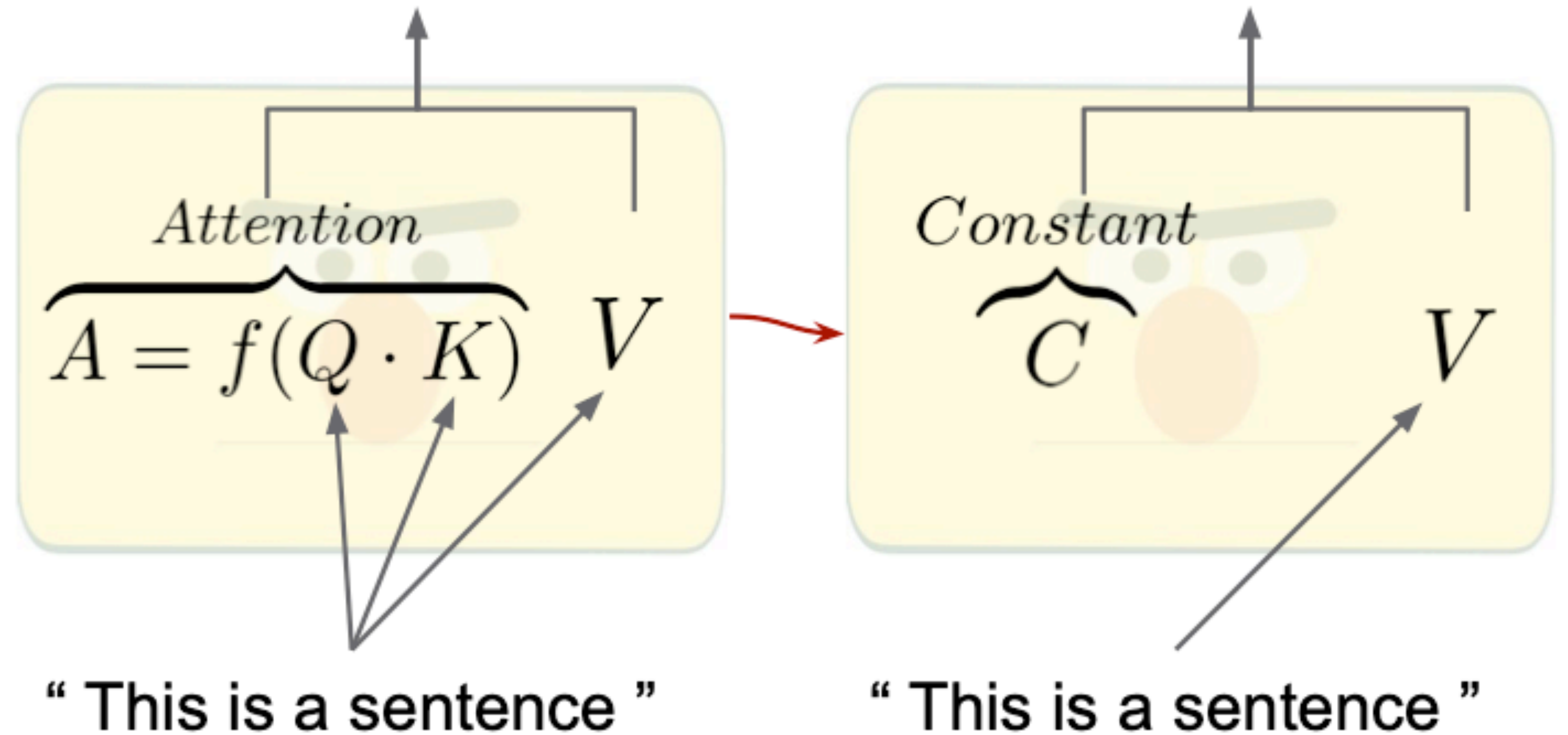


“ This is a sentence ”



How Much Does Attention Actually Attend?

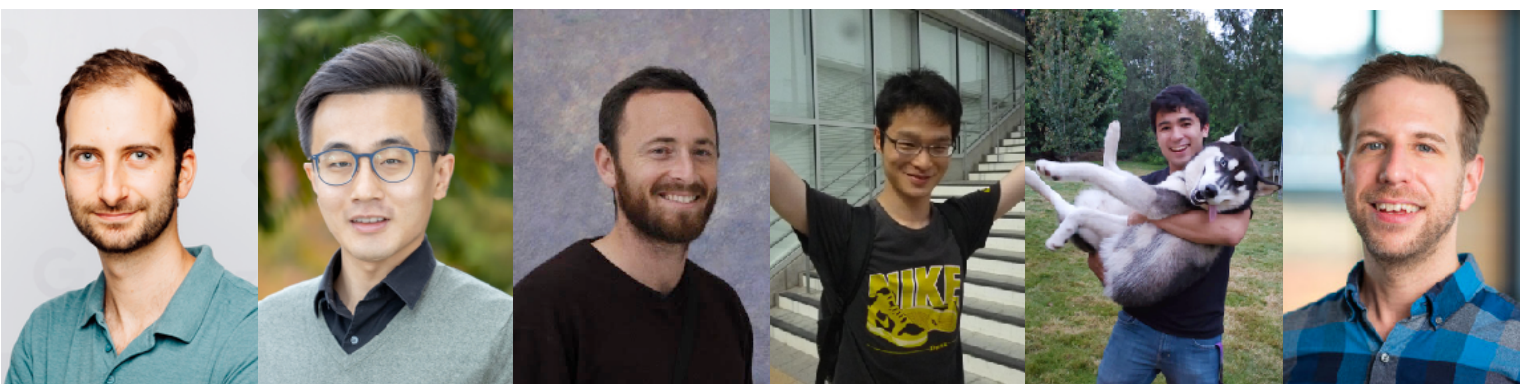
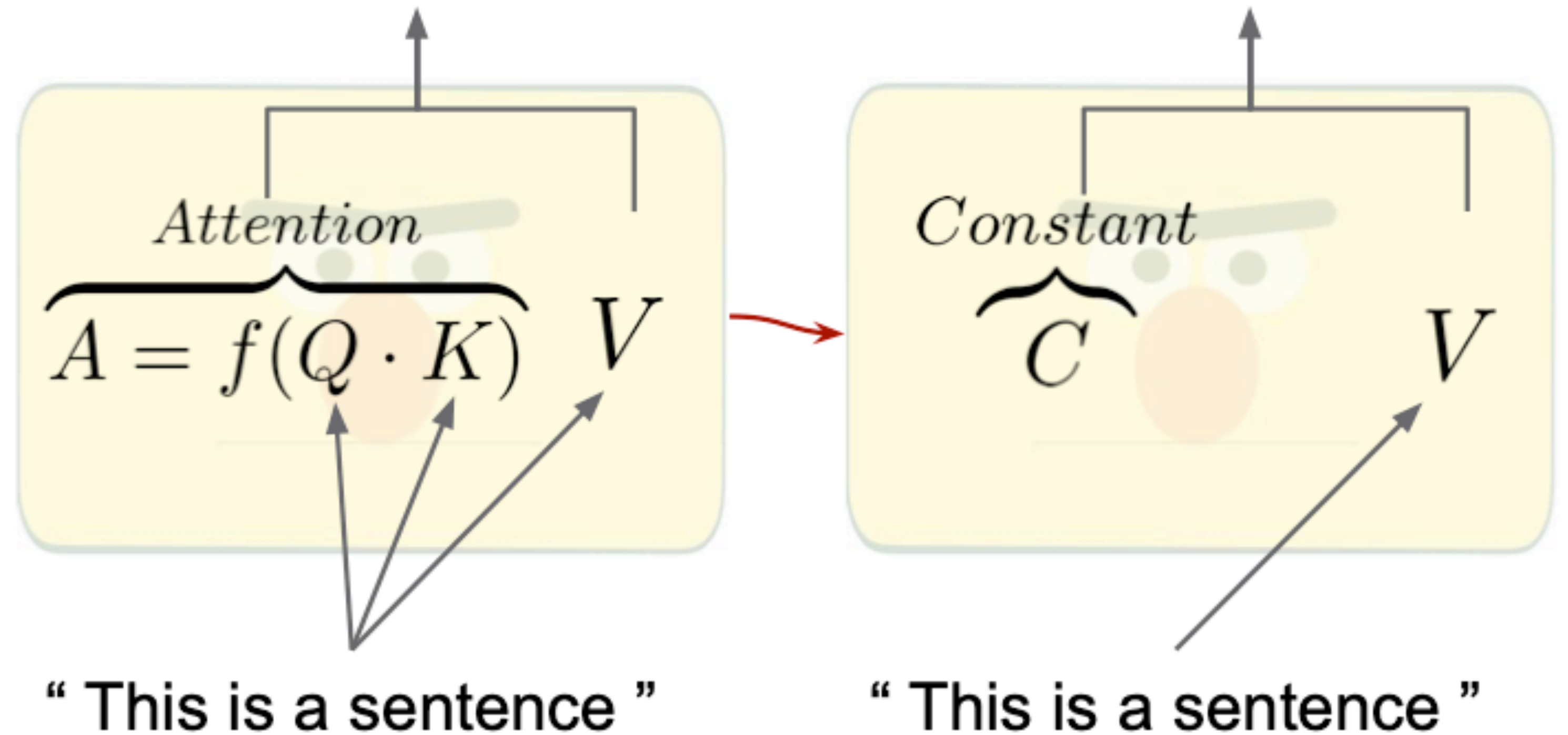
Hassid, Peng, Rotem, Kasai, Montero, Smith & S., Findings of EMNLP 2022



How Much Does Attention Actually Attend?

Hassid, Peng, Rotem, Kasai, Montero, Smith & S., Findings of EMNLP 2022

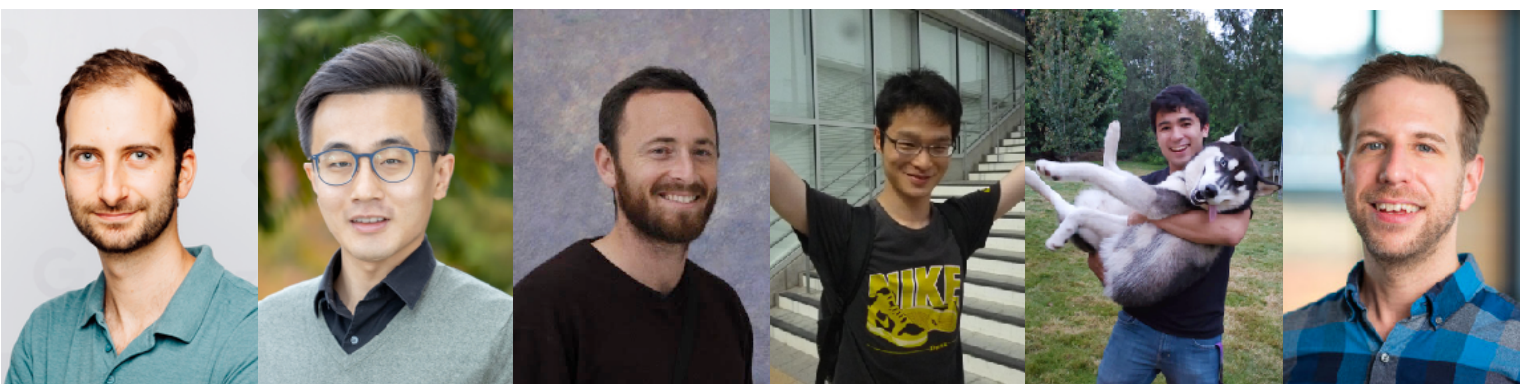
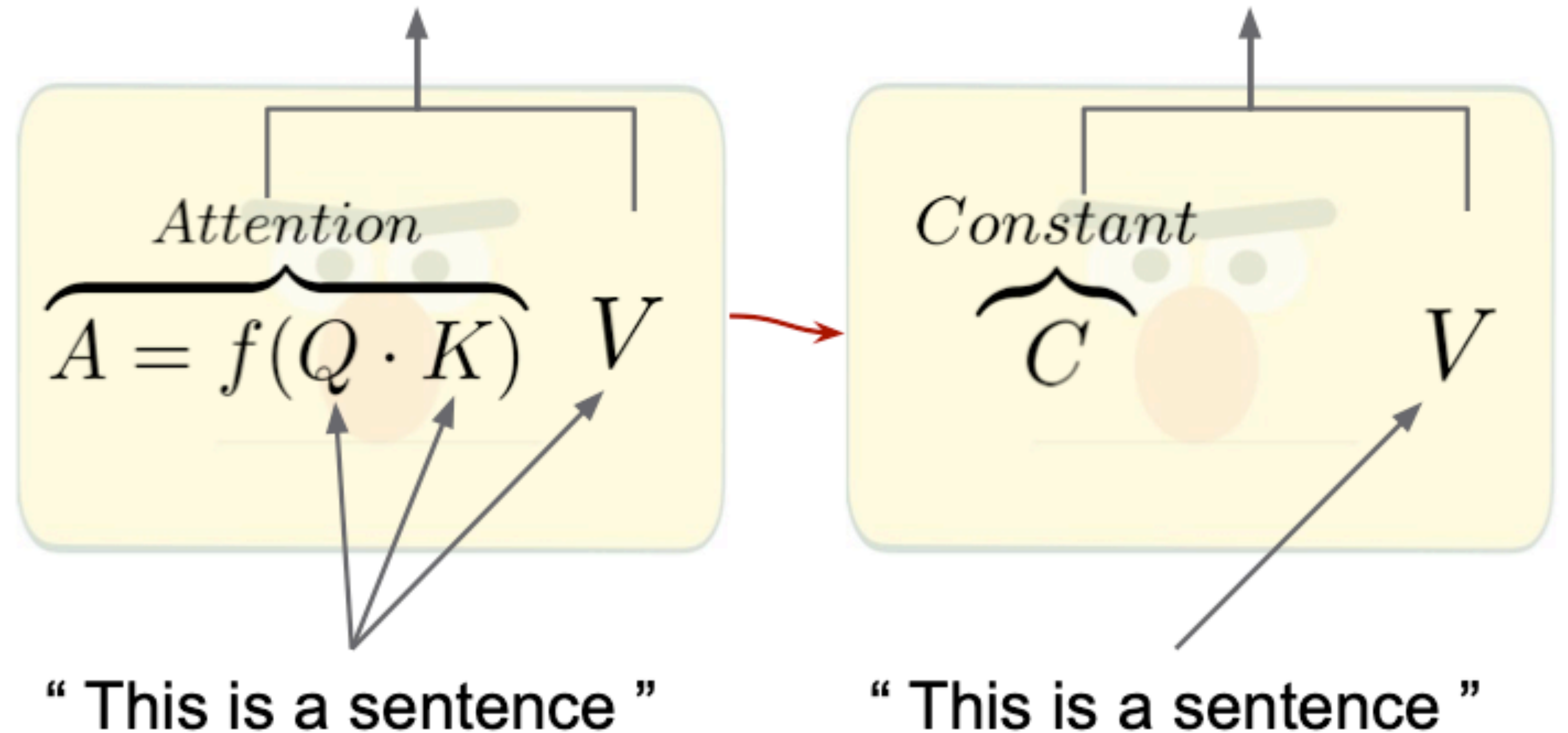
- Model doesn't collapse
 - Average accuracy loss of 8% only



How Much Does Attention Actually Attend?

Hassid, Peng, Rotem, Kasai, Montero, Smith & S., Findings of EMNLP 2022

- Model doesn't collapse
 - Average accuracy loss of 8% only
- Potential for huge savings



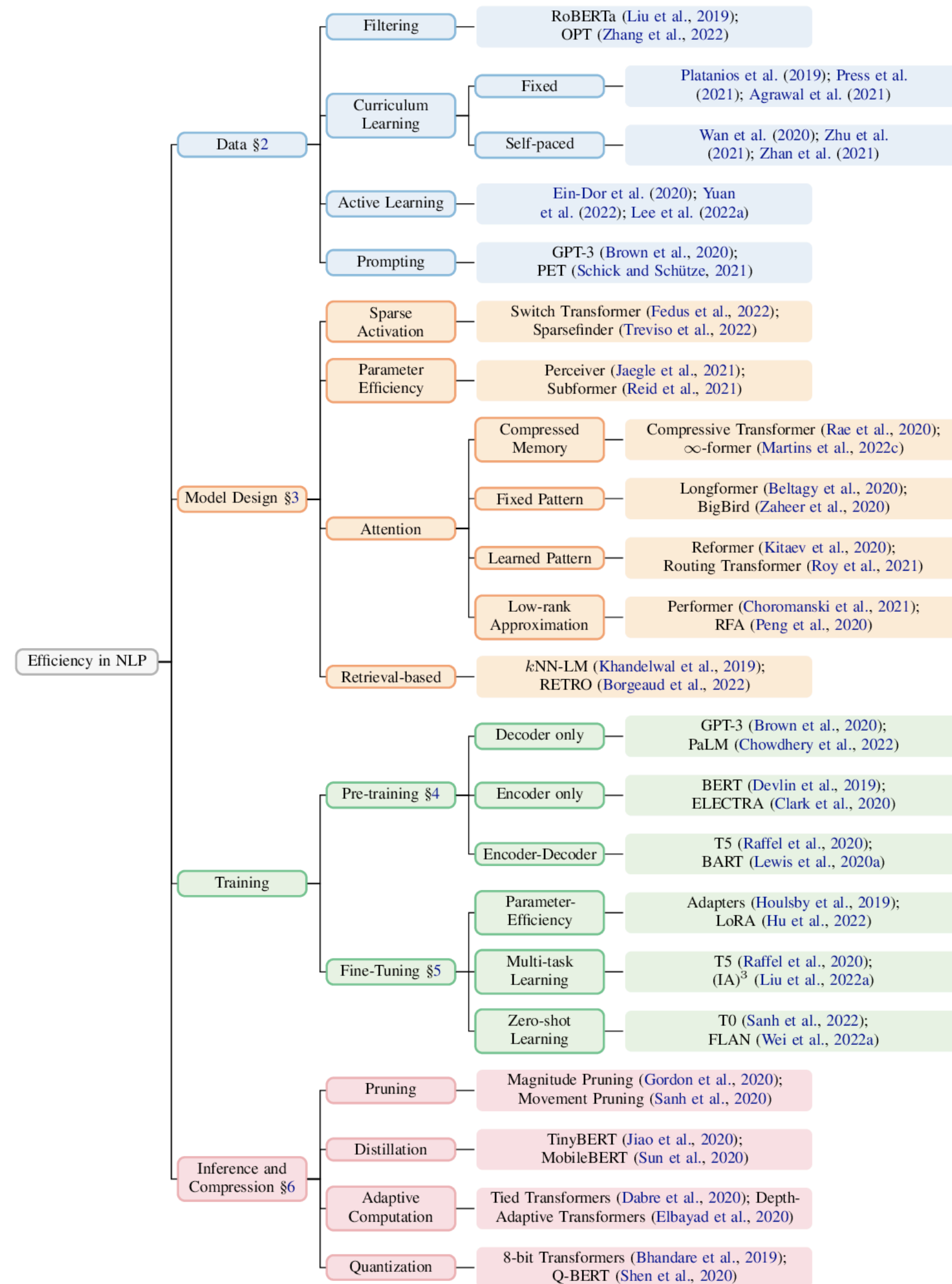
Efficient Modeling

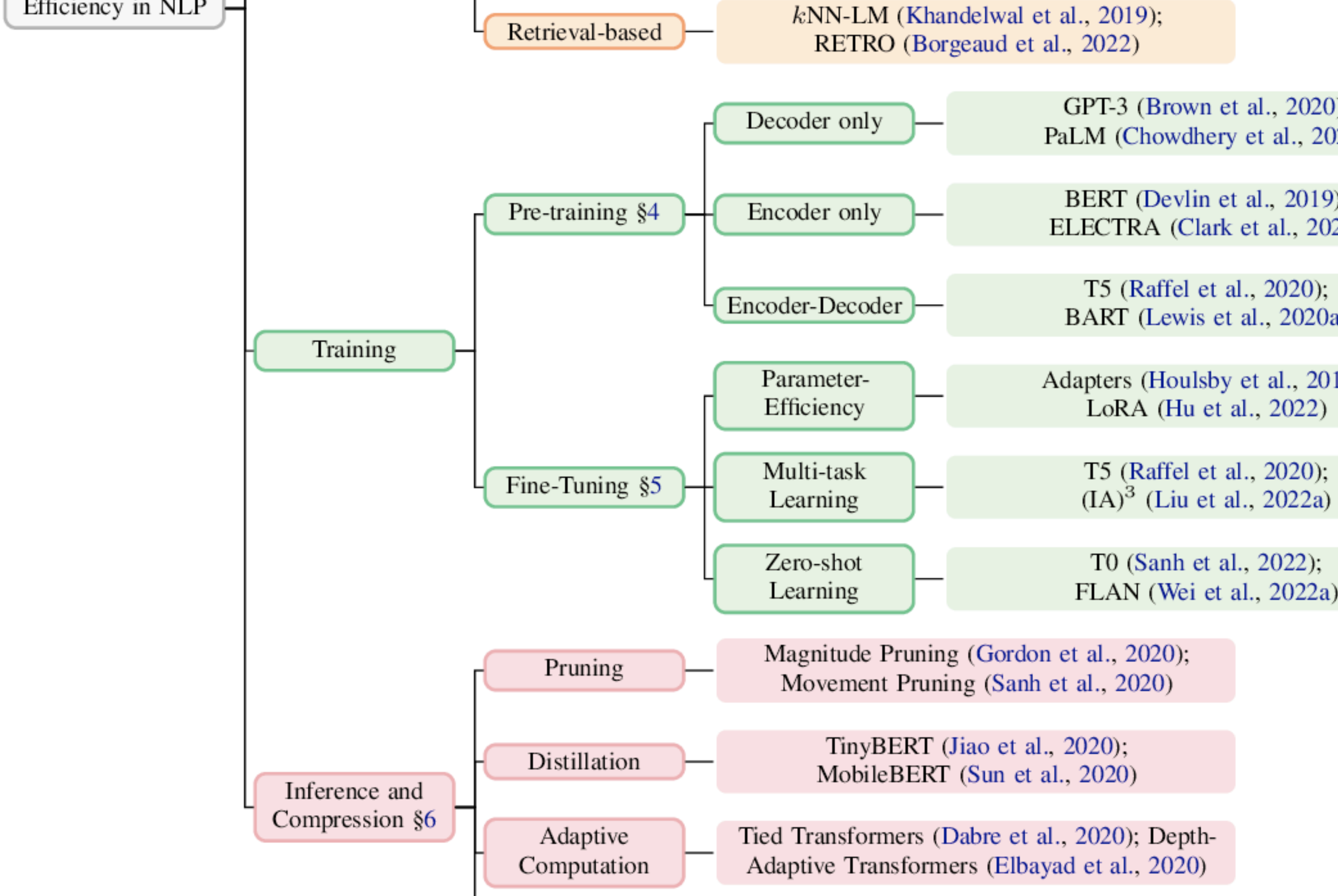
Open Questions



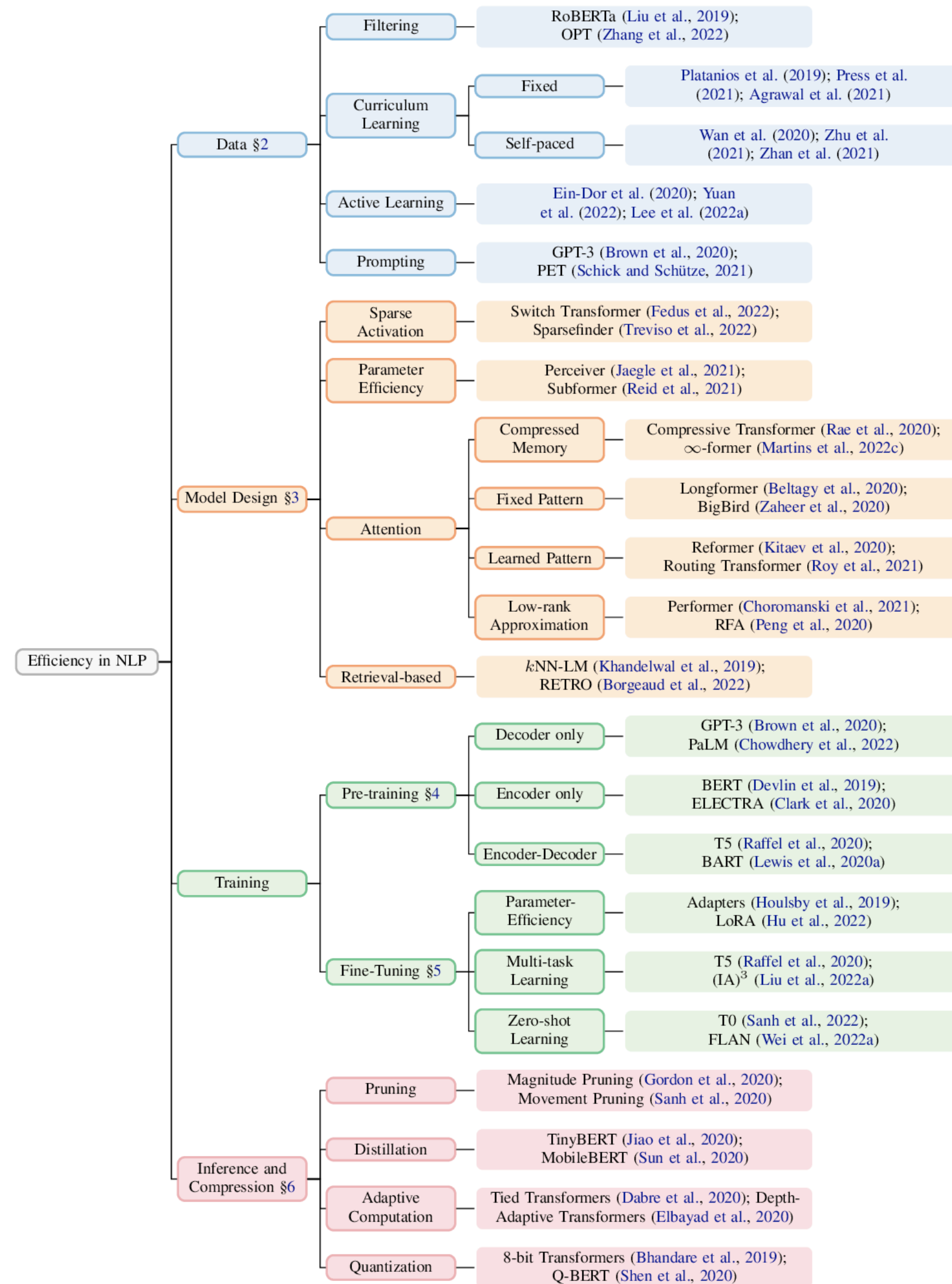
- Can we find the next generation of Transformers?
 - S4 (Gu et al., 2021)
- Should we store knowledge in the model parameters?
 - Retrieval-based models
 - Gu et al (2018); Lewis et al. (2020); Li et al. (2022); Borgeaud et al. (2022)

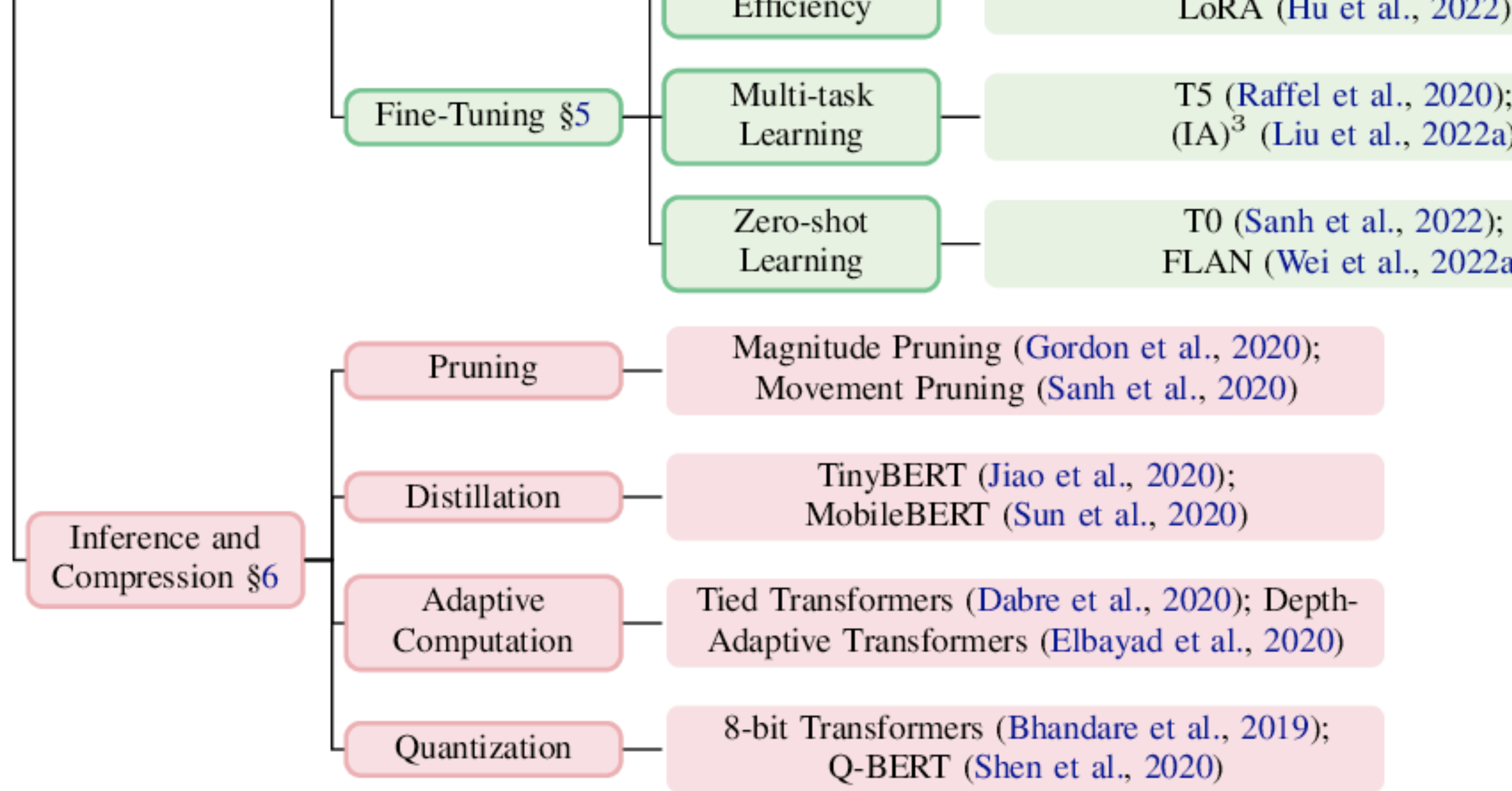
Efficient Methods in NLP



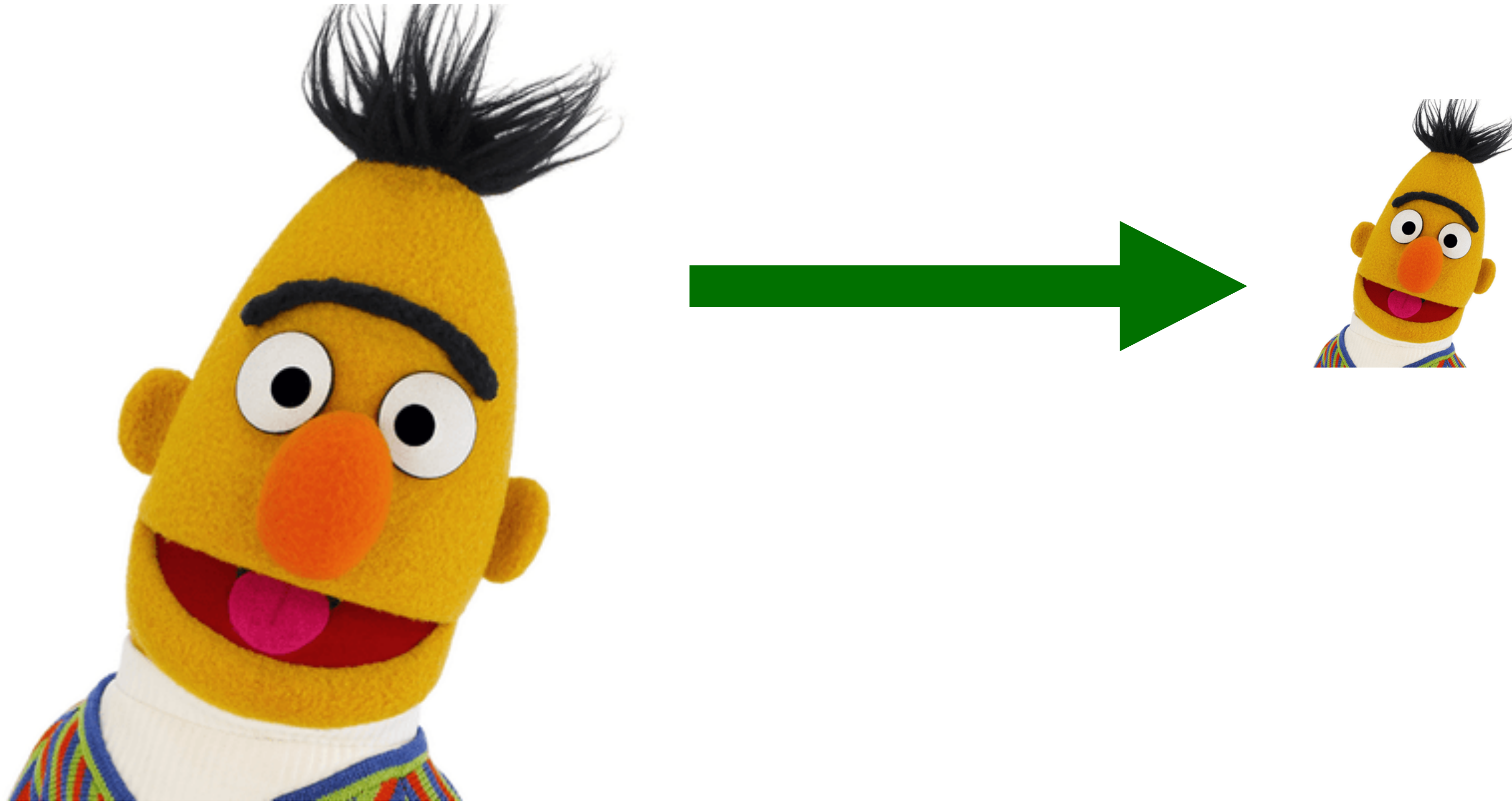


Efficient Methods in NLP





Efficient Inference



Efficient Inference



- **Model distillation**

- Aka, student/teacher model
- Hinton et al., 2015; Sun et al., 2019; Sanh et al., 2019

Efficient Inference



- **Model distillation**

- Aka, student/teacher model
- Hinton et al., 2015; Sun et al., 2019; Sanh et al., 2019

- **Pruning / Structural Pruning**

- Han et al., 2016; Lee et al., 2019; Frankle & Corbin, 2019; Gordon et al., 2018; Michel et al., 2019; Fan et al., 2020
- Dodge, S., et al., 2019

Efficient Inference



- **Model distillation**

- Aka, student/teacher model
- Hinton et al., 2015; Sun et al., 2019; Sanh et al., 2019

- **Pruning / Structural Pruning**

- Han et al., 2016; Lee et al., 2019; Frankle & Corbin, 2019; Gordon et al., 2018; Michel et al., 2019; Fan et al., 2020
- Dodge, S., et al., 2019

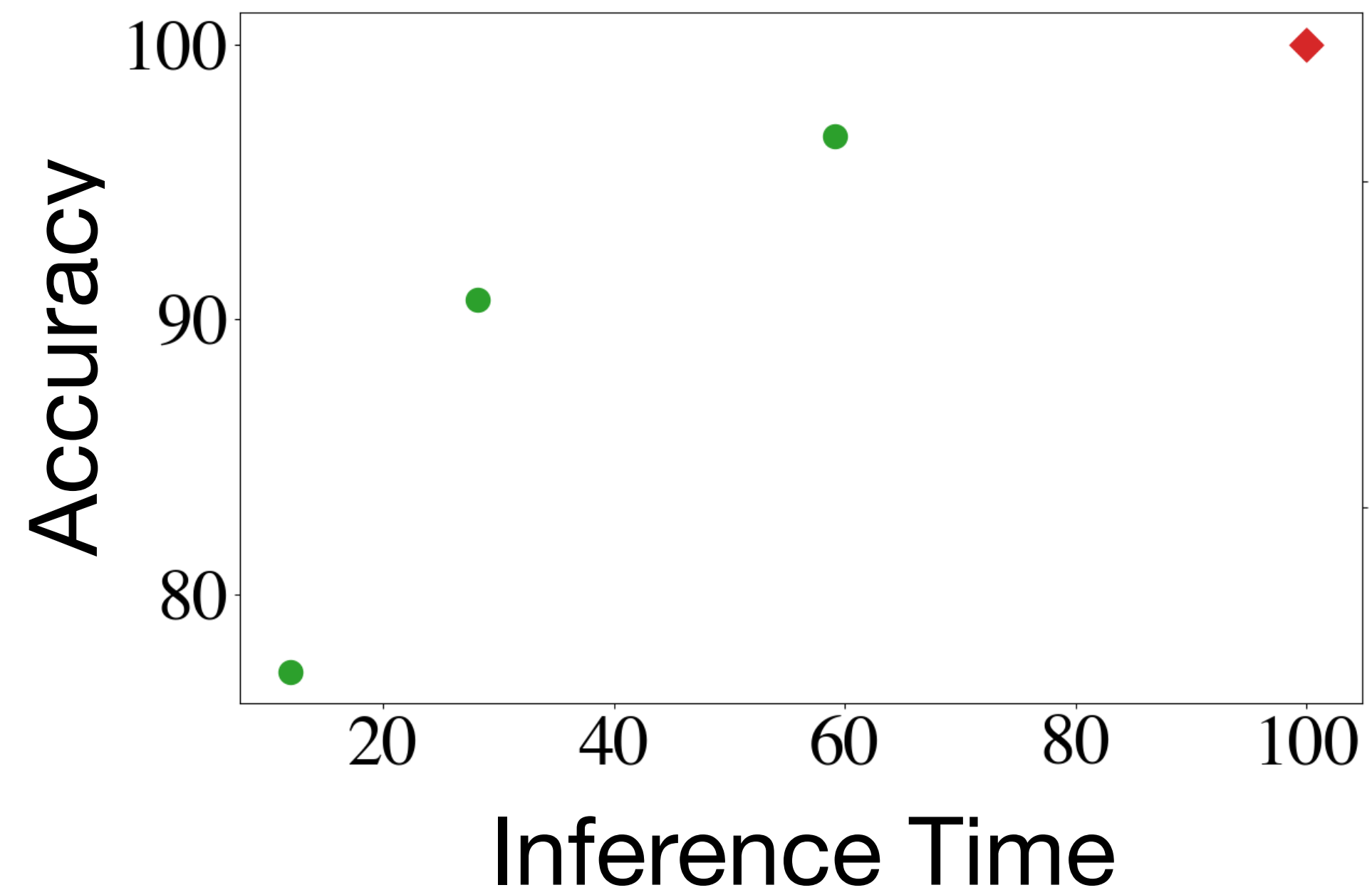
- **Quantization**

- Gong et al., 2014; Zafrir et al., 2019; Shen et al., 2019

Matching Model and Instance Complexity

S. et al., ACL 2020

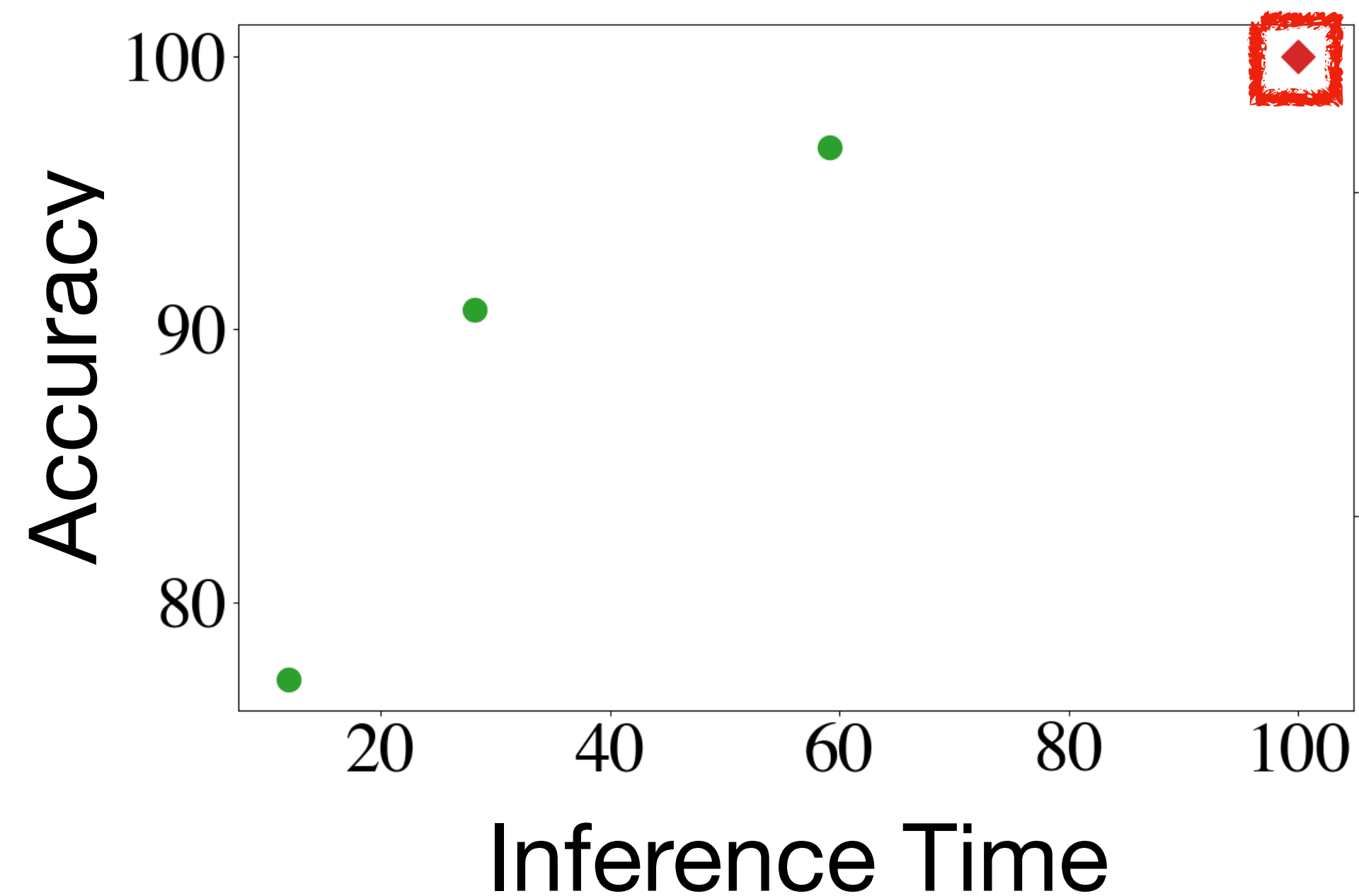
Run an *efficient* model on “*easy*” instances,
and an *expensive* model on “*hard*” instances



Matching Model and Instance Complexity

S. et al., ACL 2020

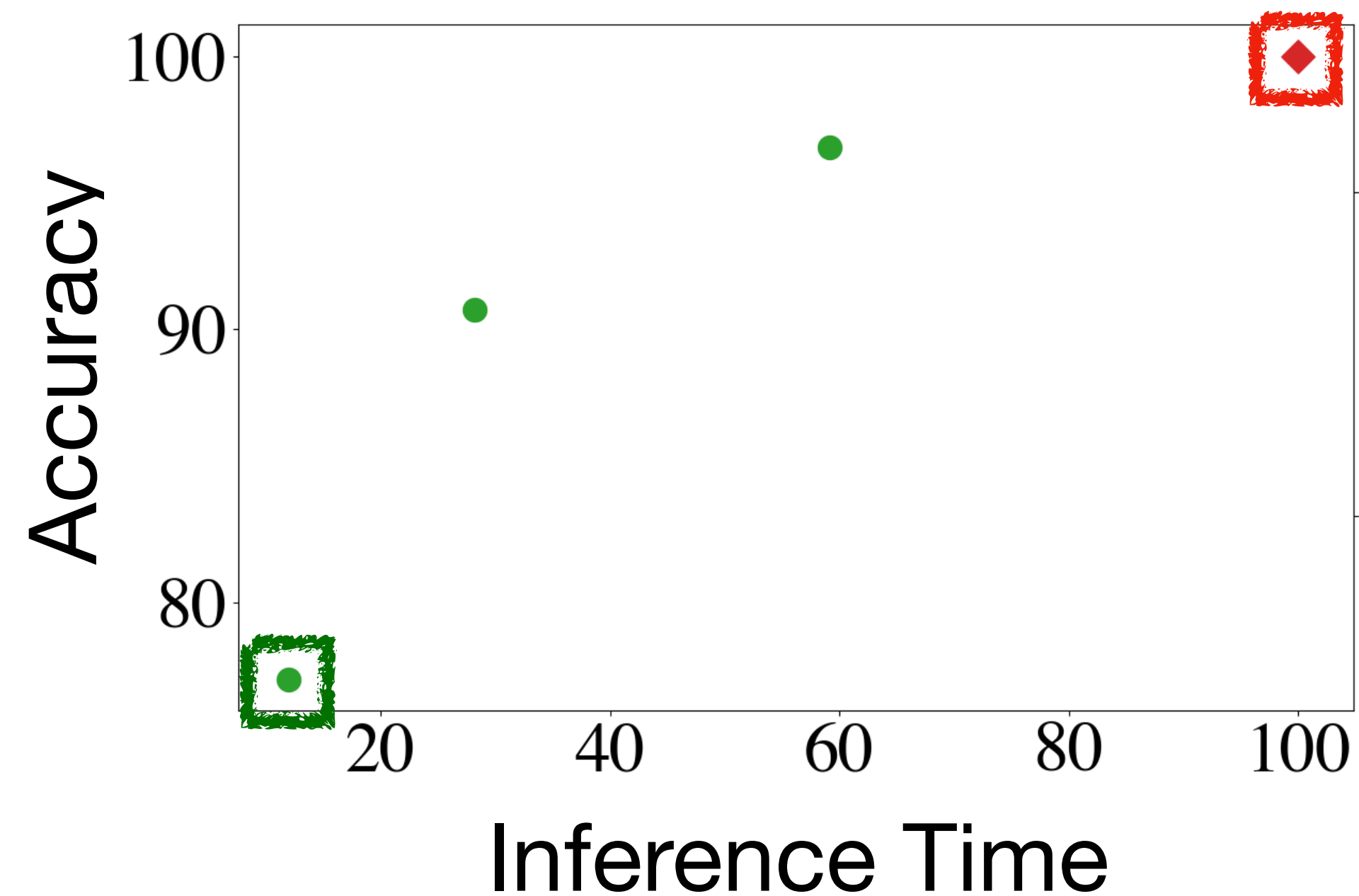
Run an *efficient* model on “*easy*” instances,
and an *expensive* model on “*hard*” instances



Matching Model and Instance Complexity

S. et al., ACL 2020

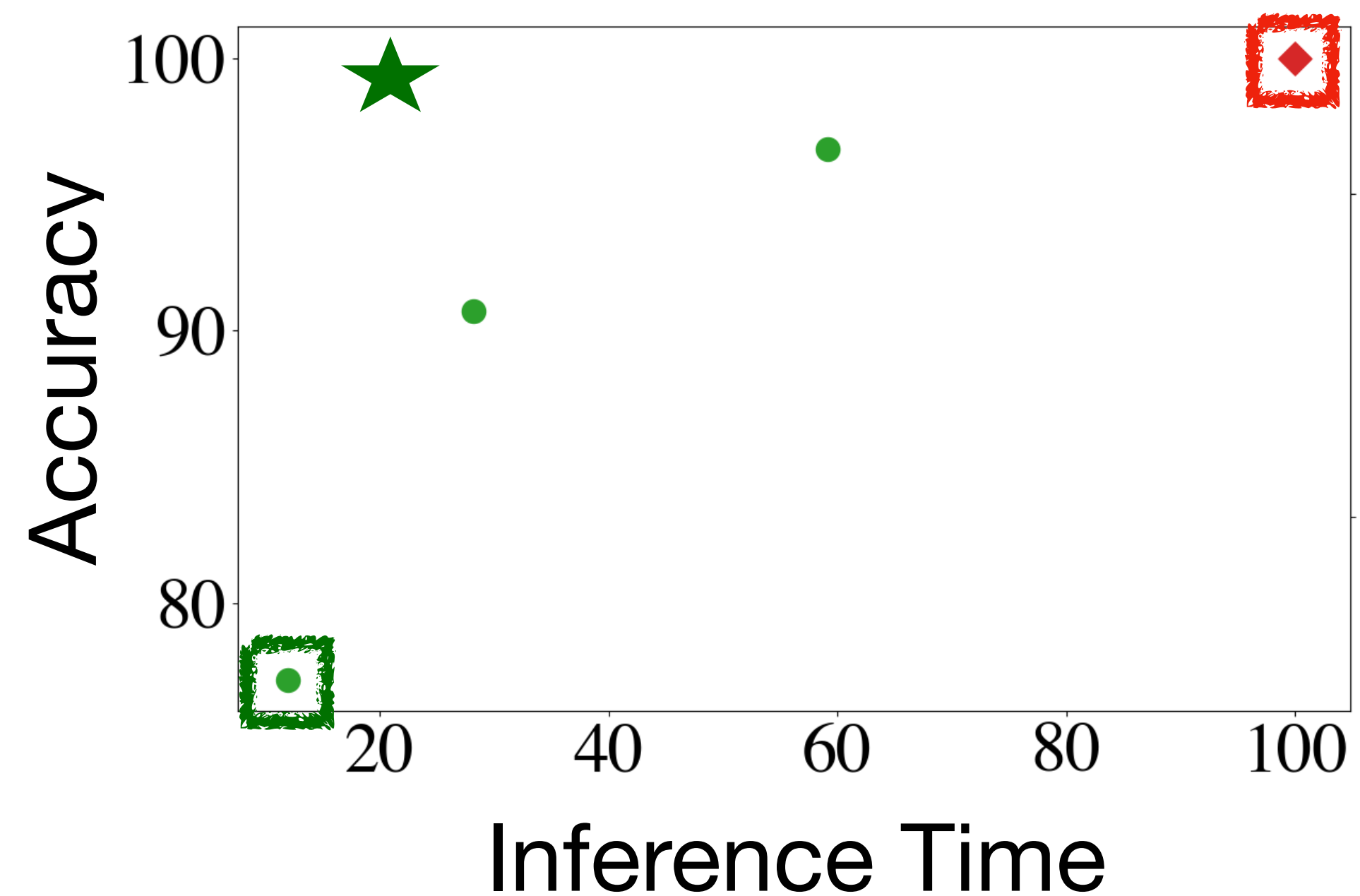
Run an *efficient* model on “*easy*” instances,
and an *expensive* model on “*hard*” instances



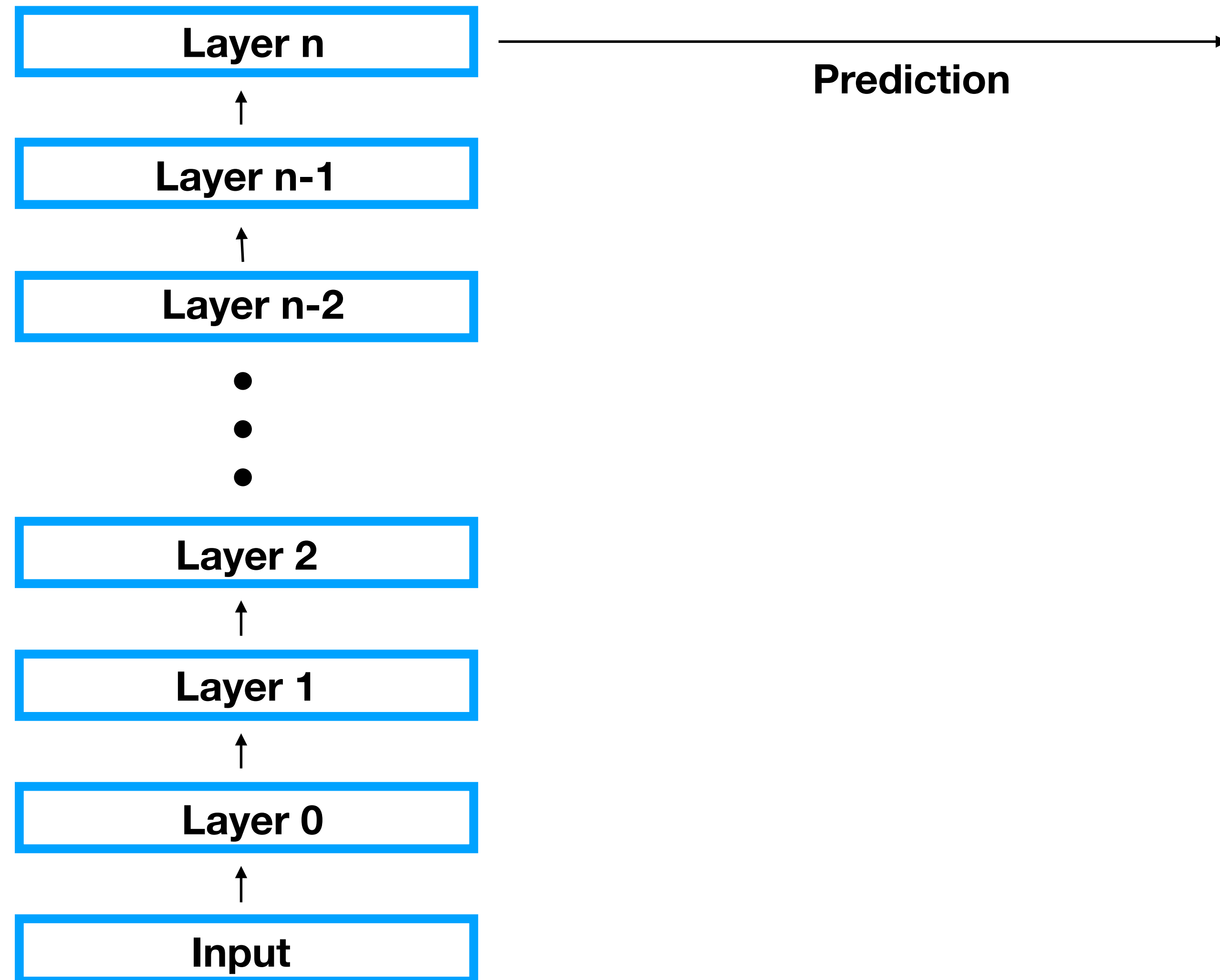
Matching Model and Instance Complexity

S. et al., ACL 2020

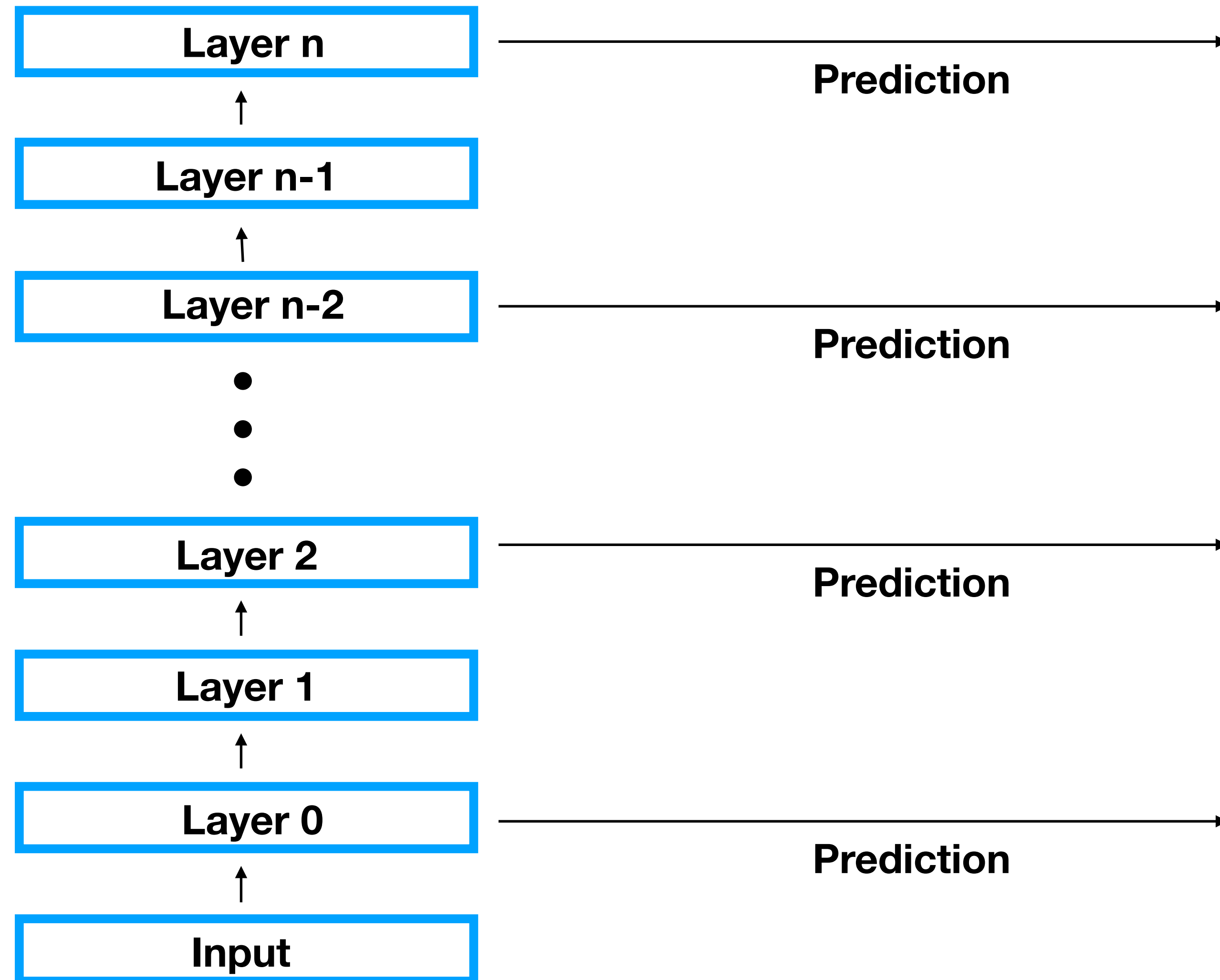
Run an *efficient* model on “*easy*” instances,
and an *expensive* model on “*hard*” instances



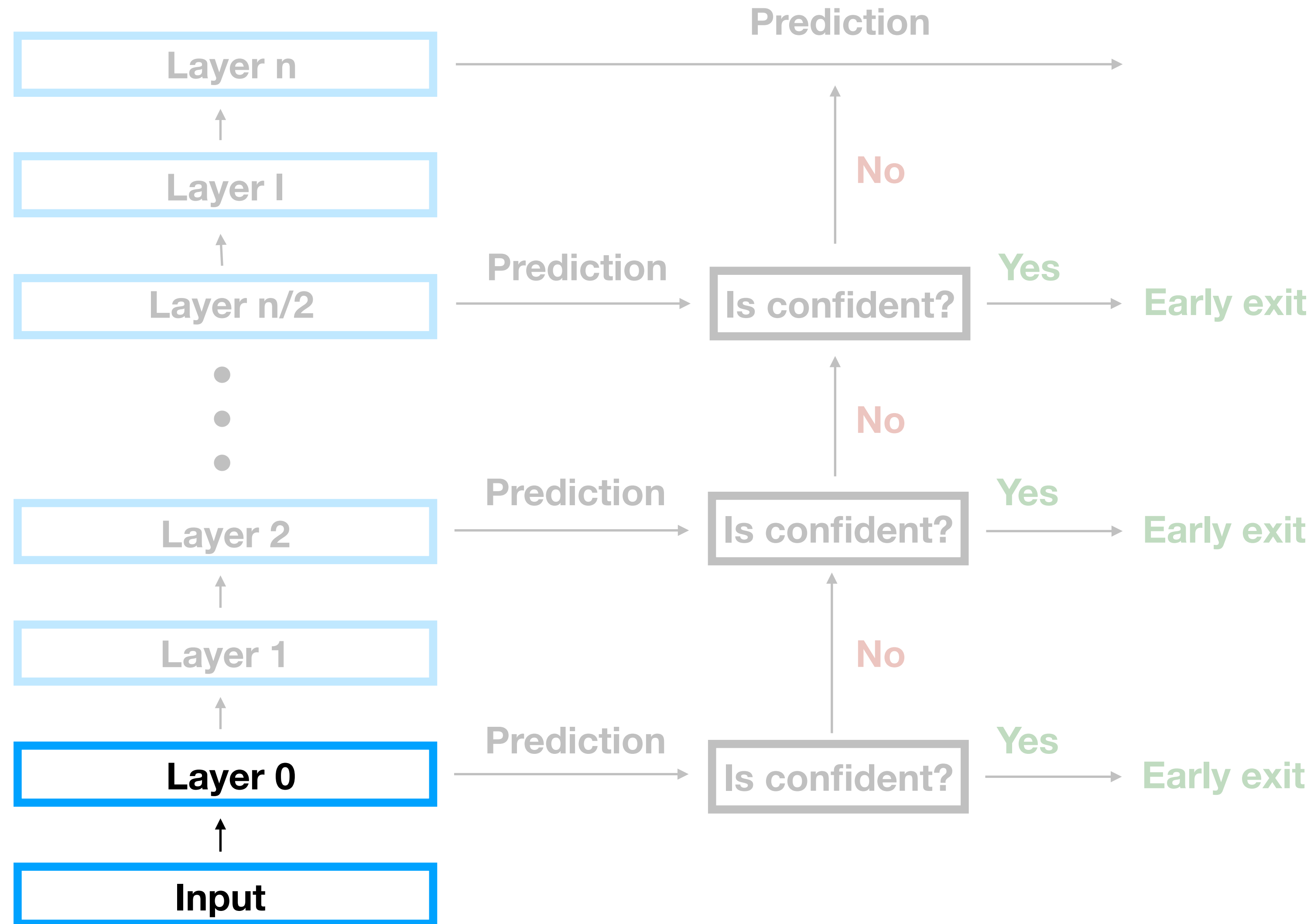
Our Approach: Training Time



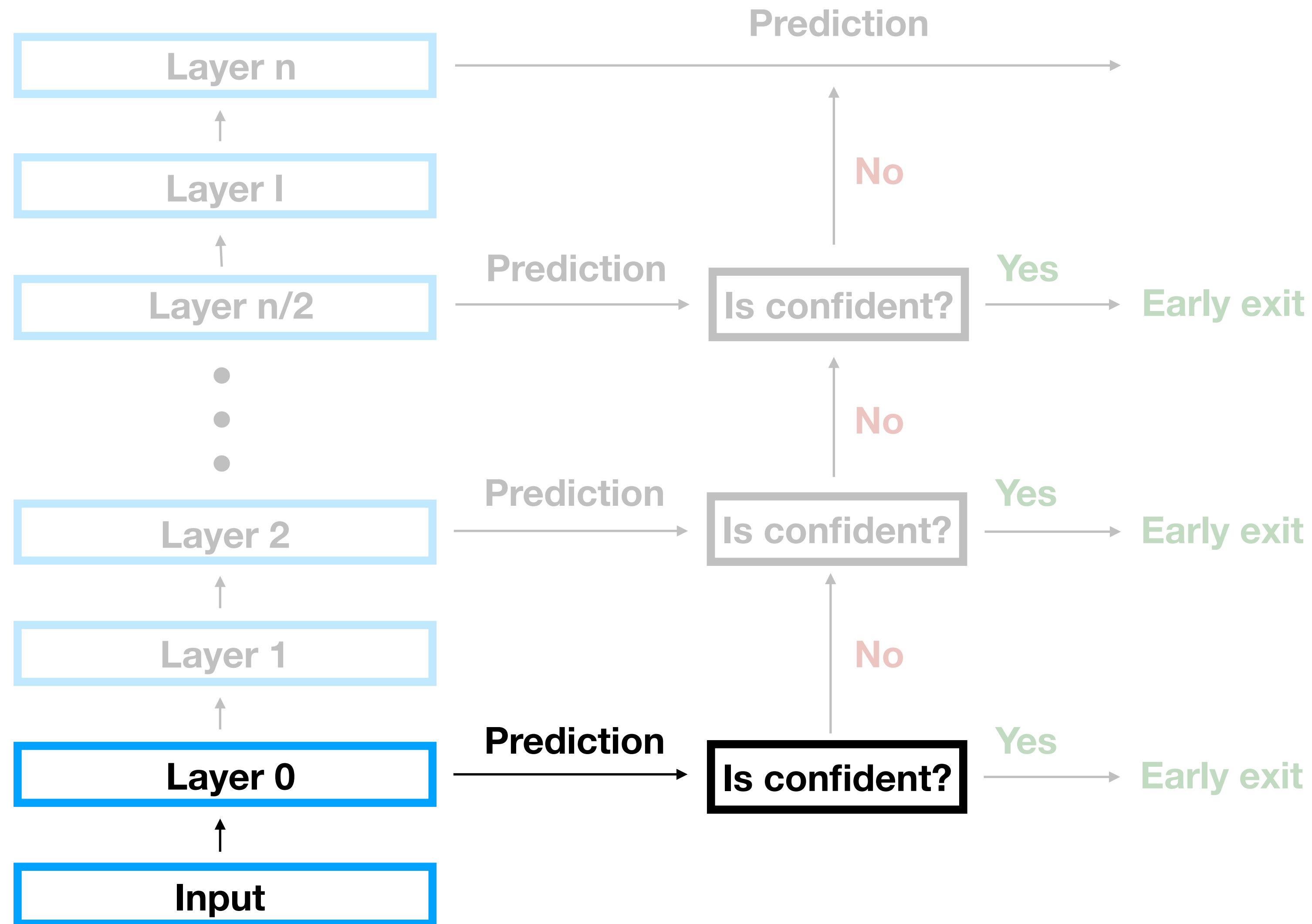
Our Approach: Training Time



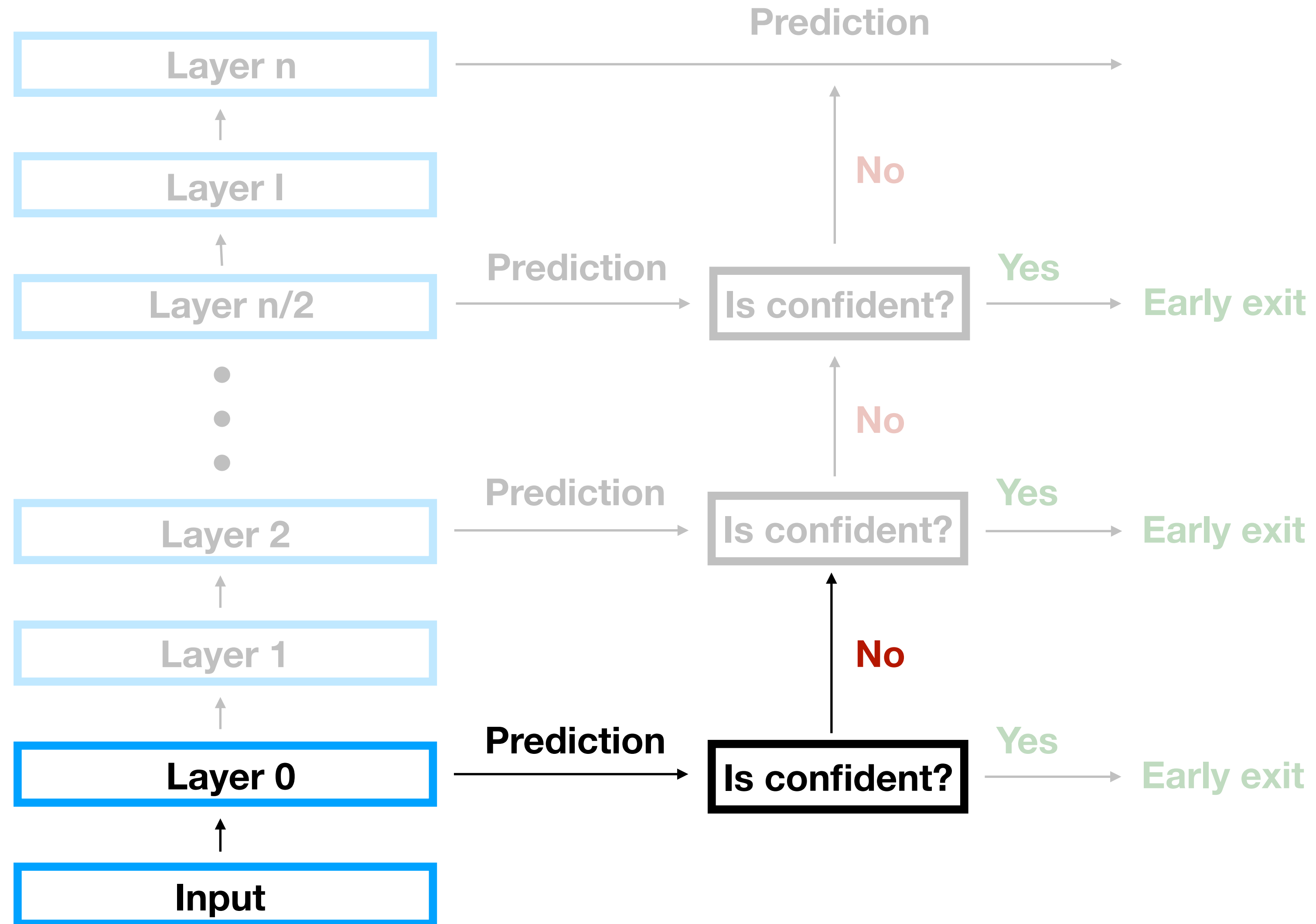
Our Approach: Test Time



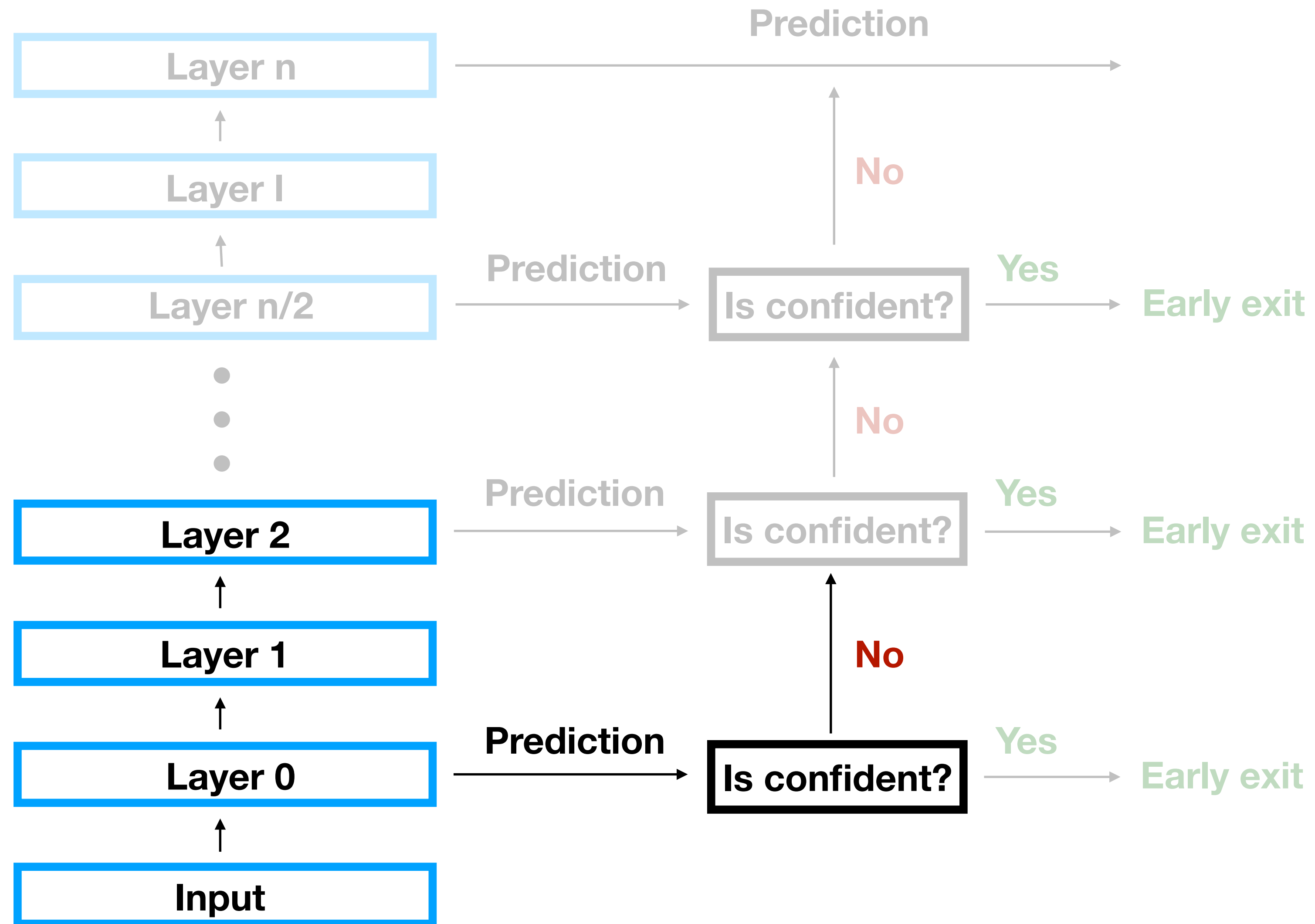
Our Approach: Test Time



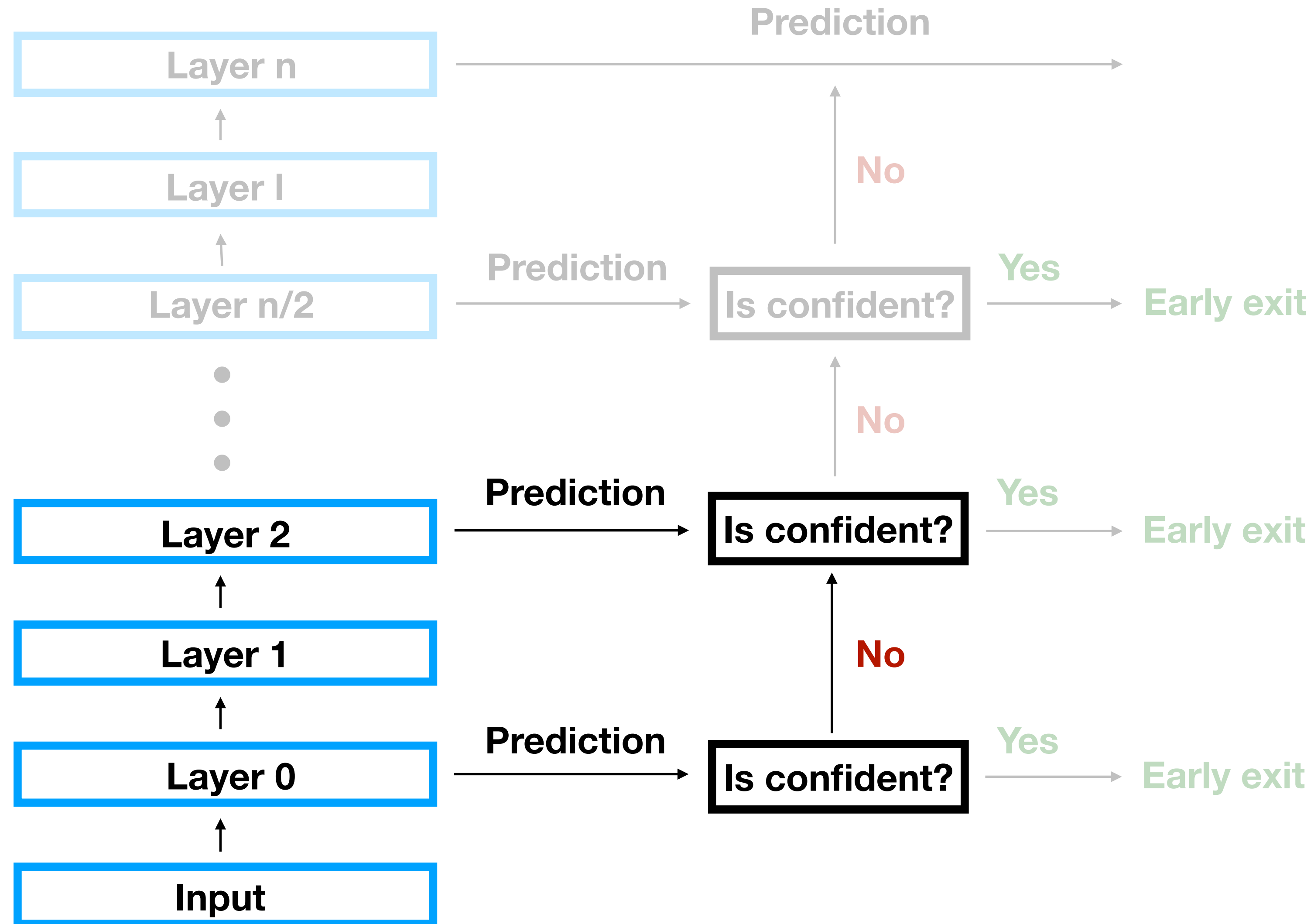
Our Approach: Test Time



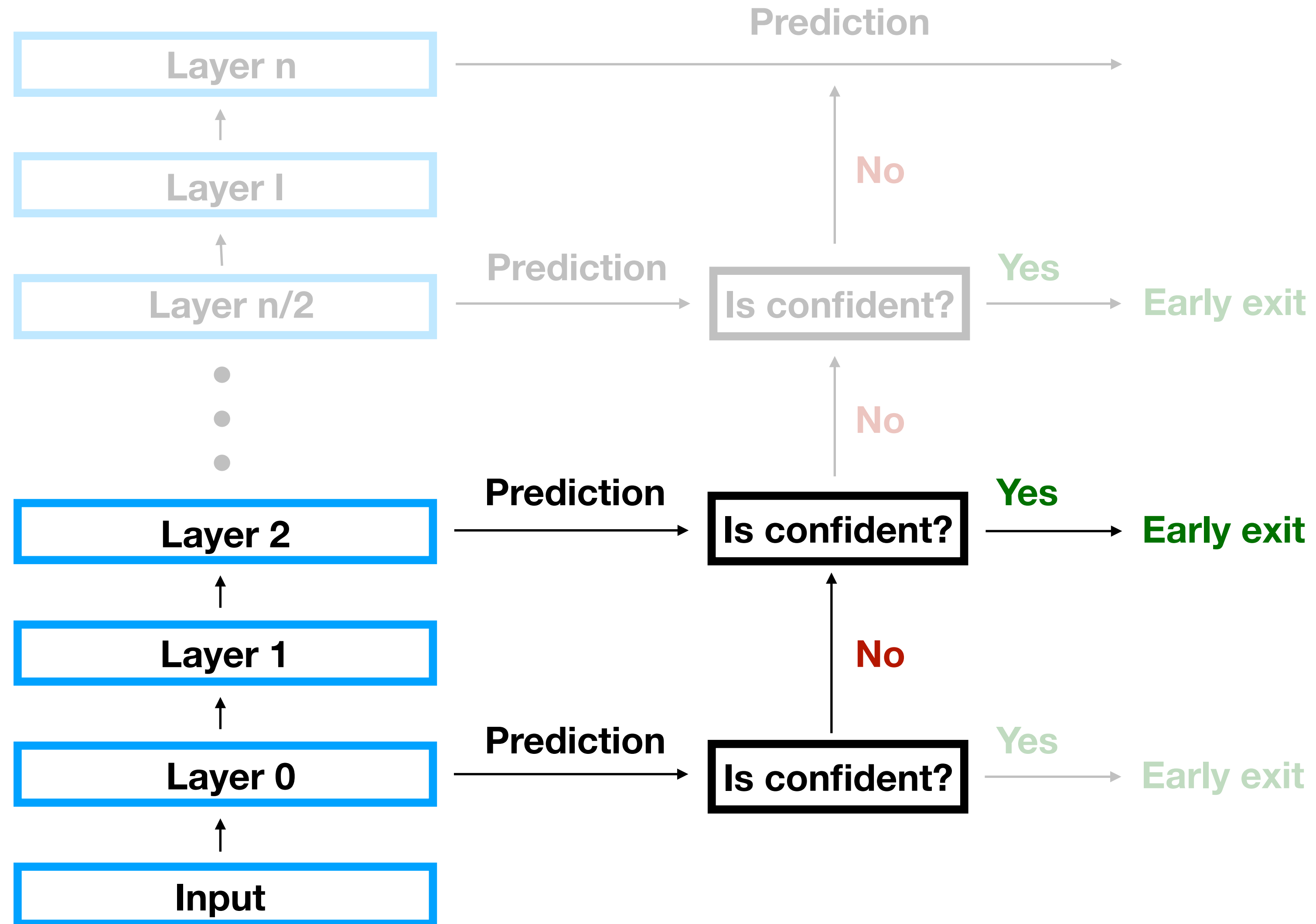
Our Approach: Test Time



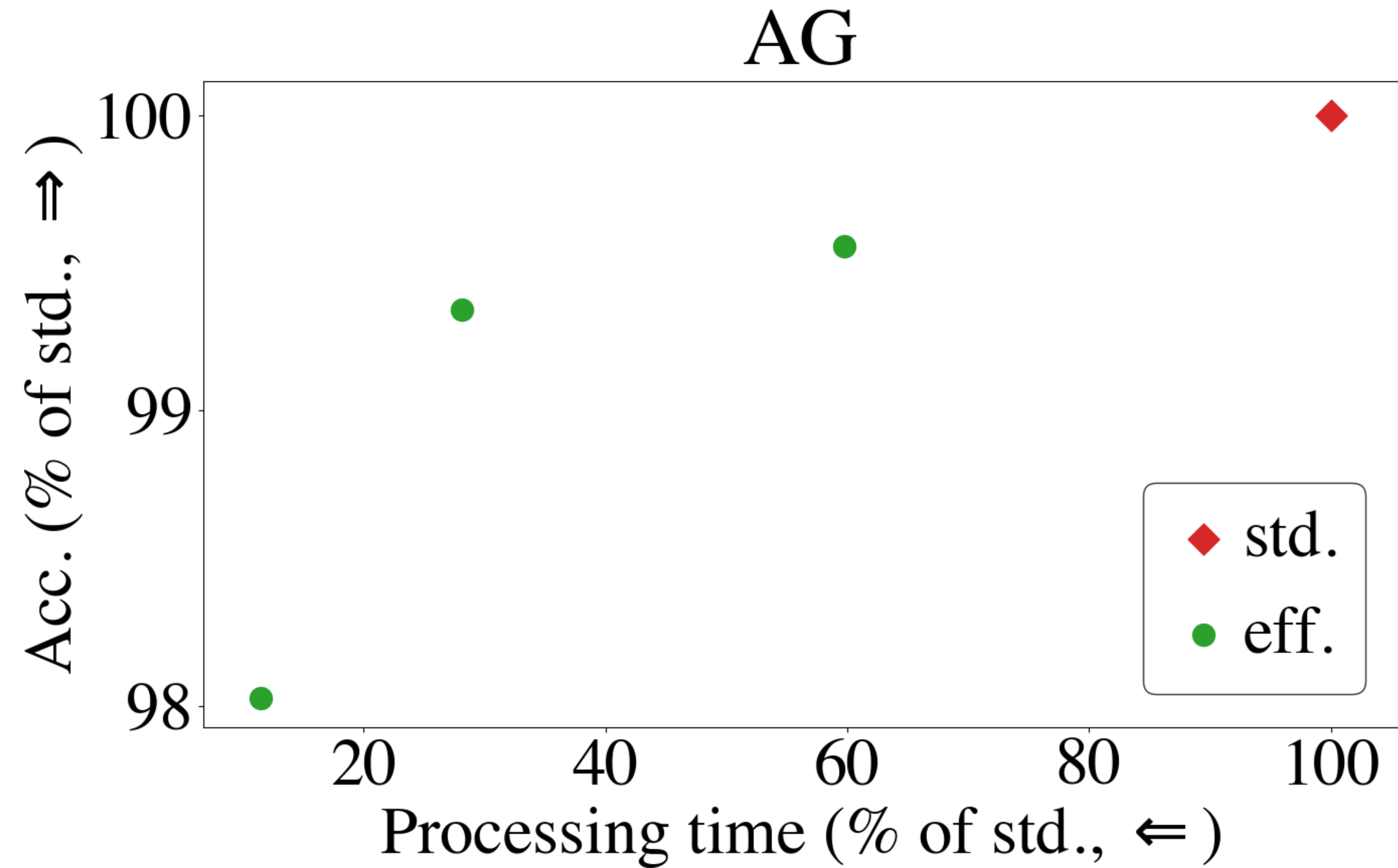
Our Approach: Test Time



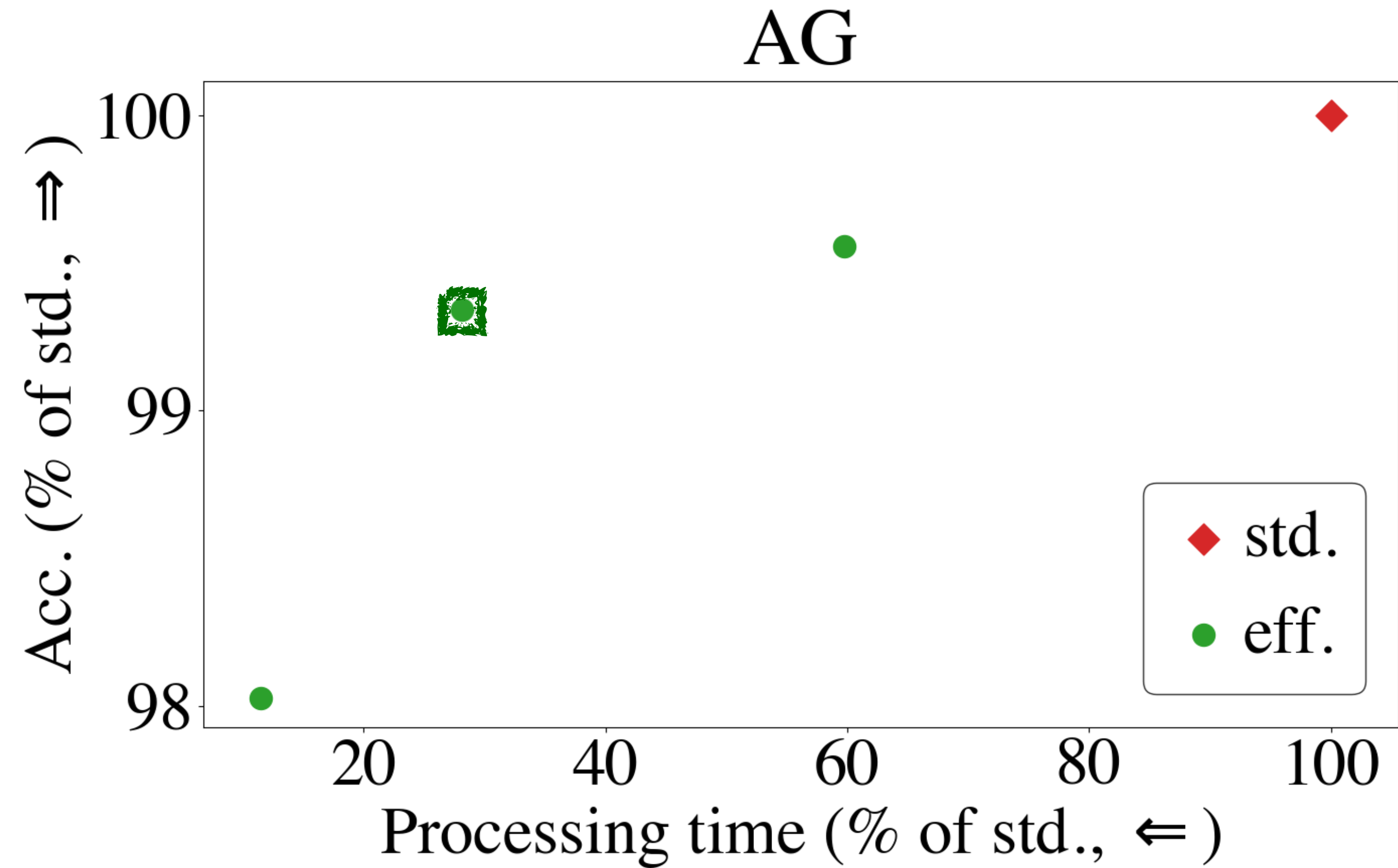
Our Approach: Test Time



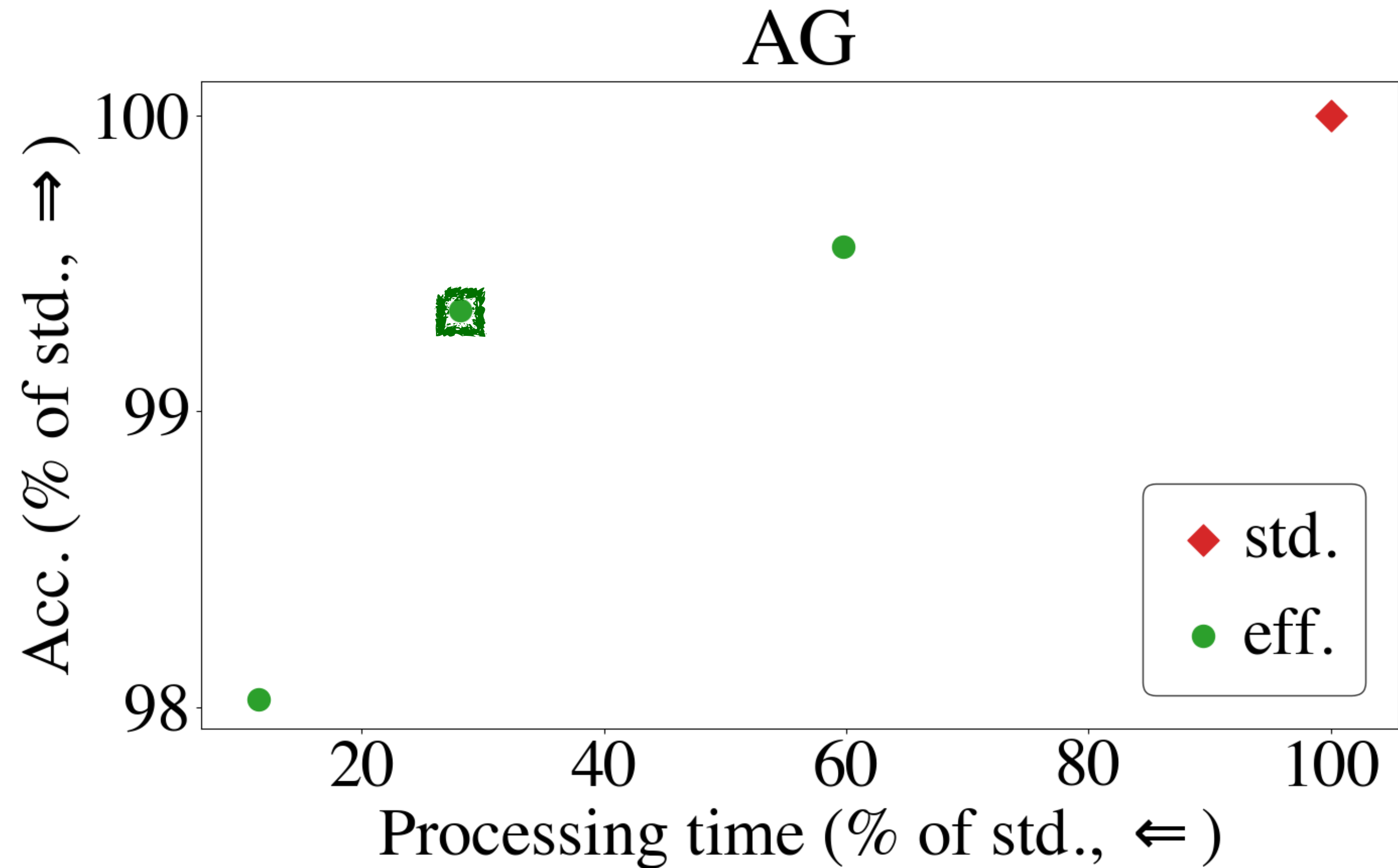
Strong Baselines!



Strong Baselines!



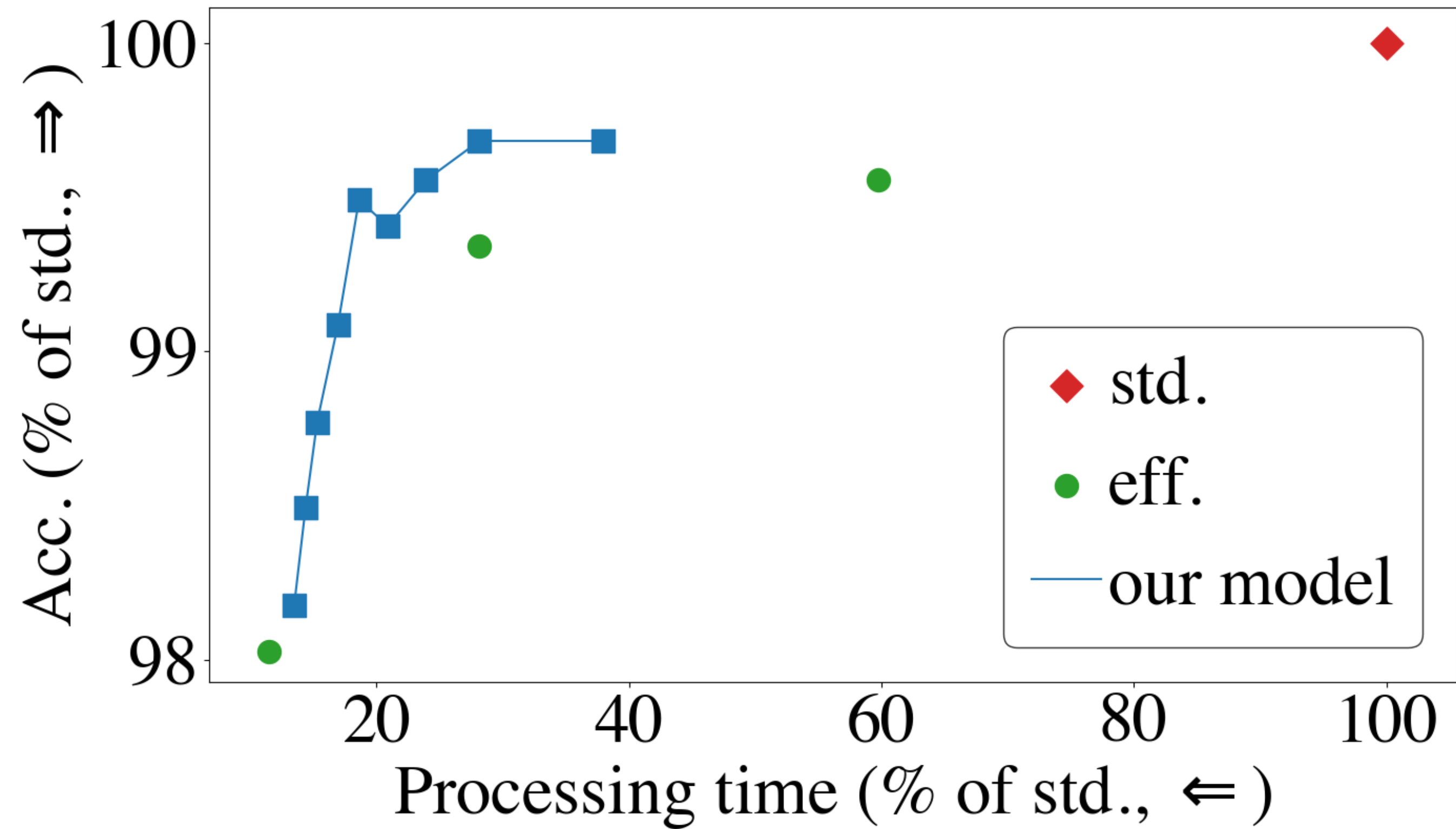
Strong Baselines!



3 times faster, within 1% of full model

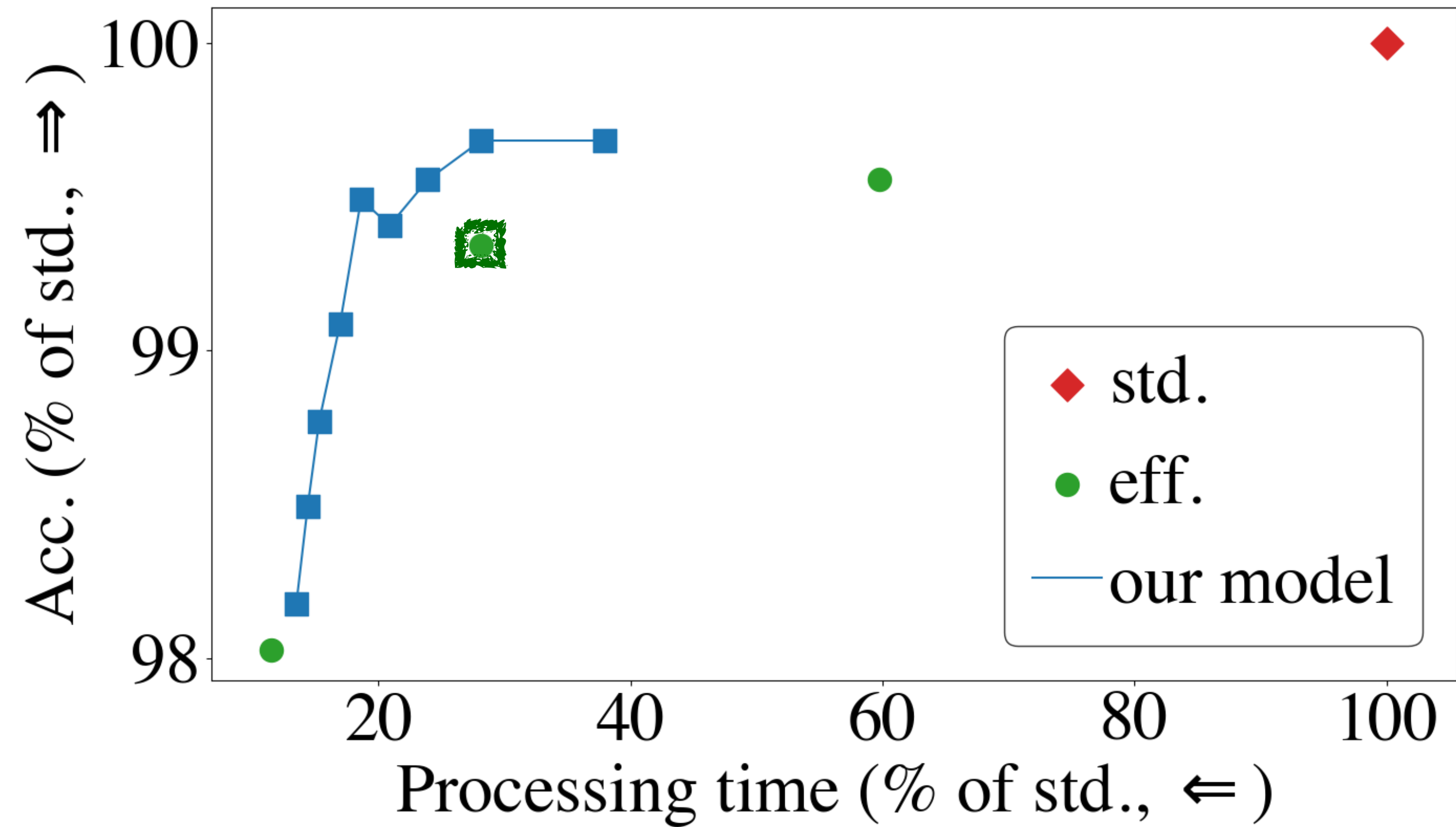
Better Speed/Accuracy Tradeoff

AG



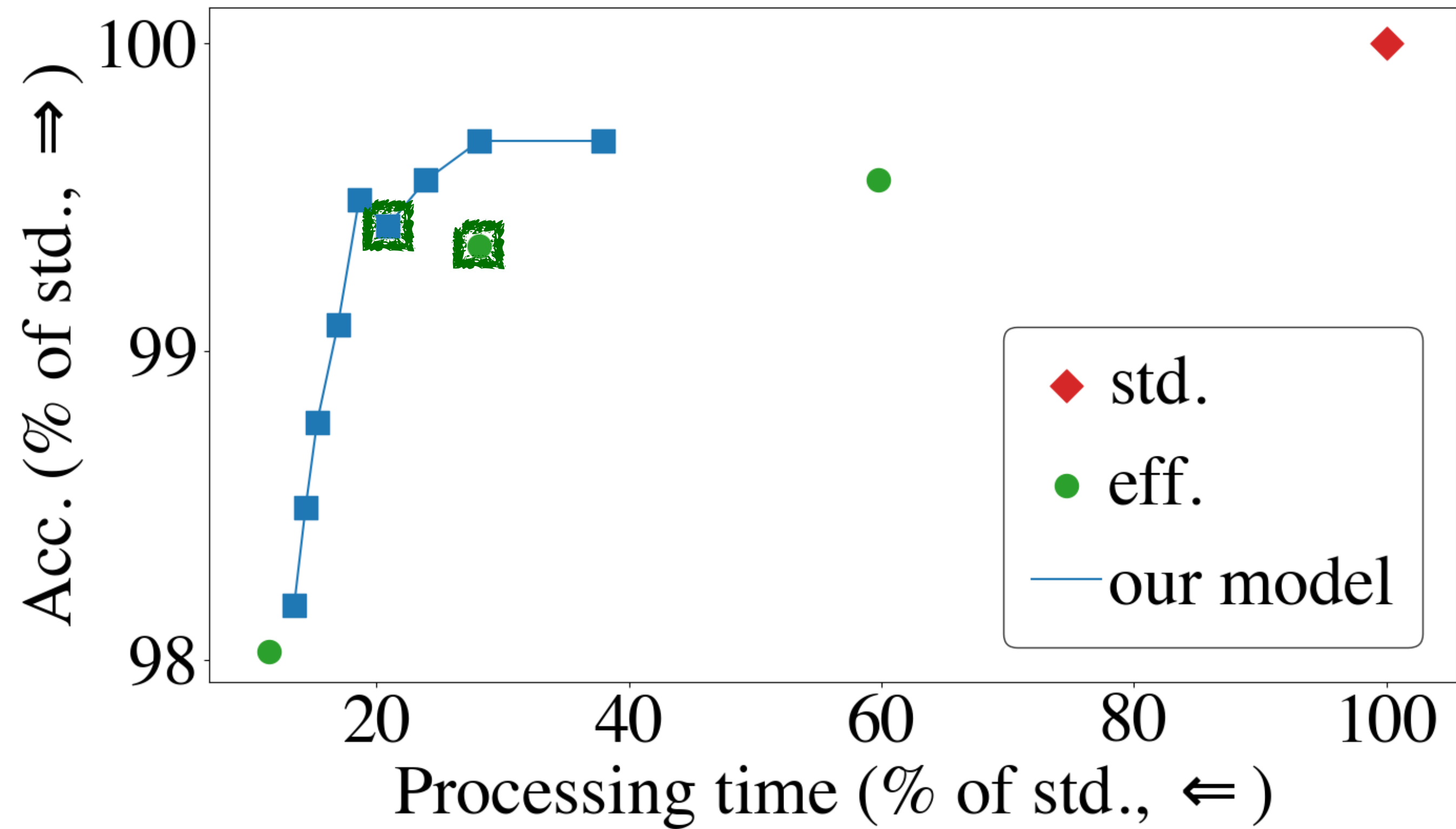
Better Speed/Accuracy Tradeoff

AG



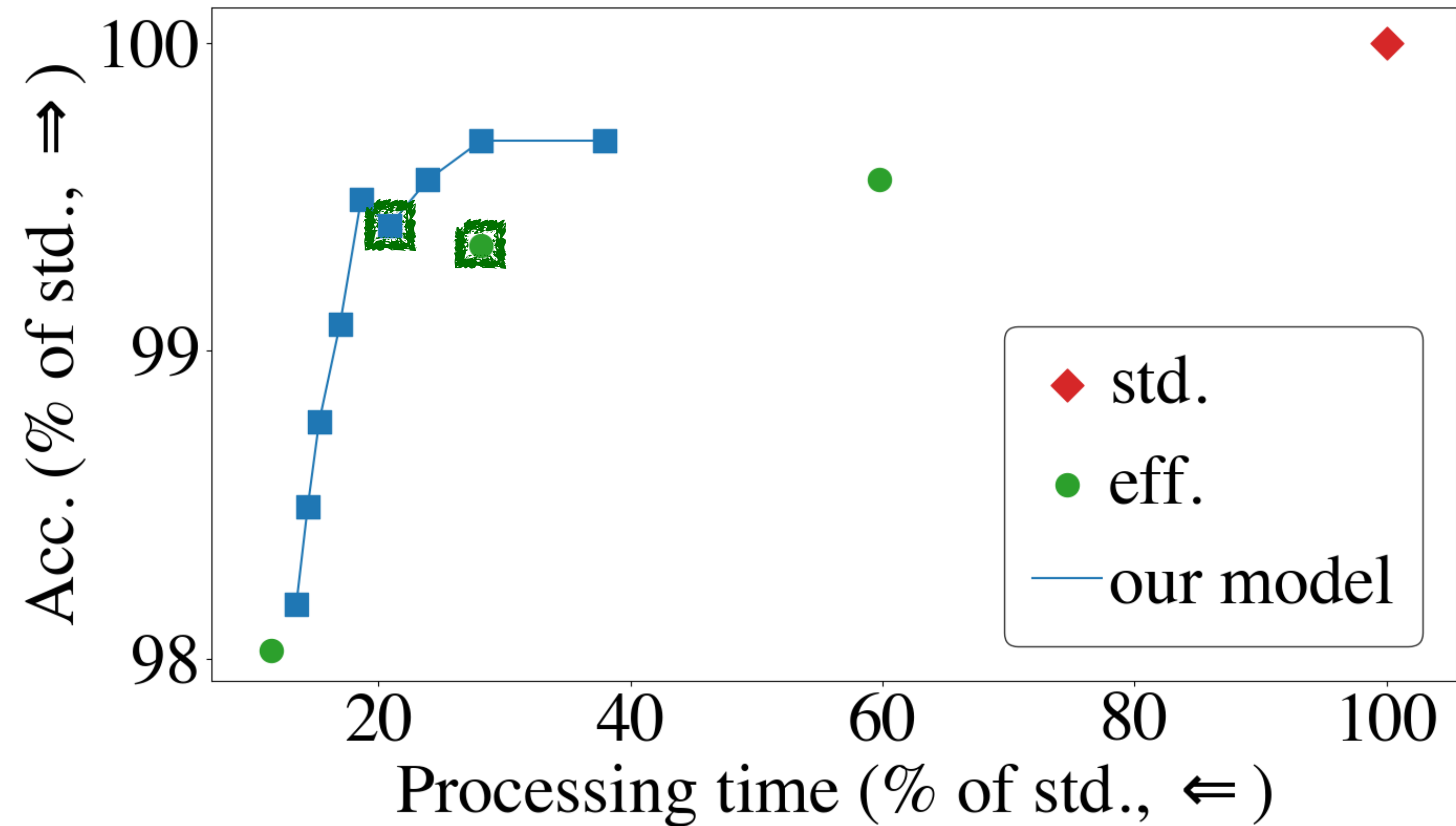
Better Speed/Accuracy Tradeoff

AG



Better Speed/Accuracy Tradeoff

AG



5 times faster, within 1% of full model



Efficiency

Open Questions

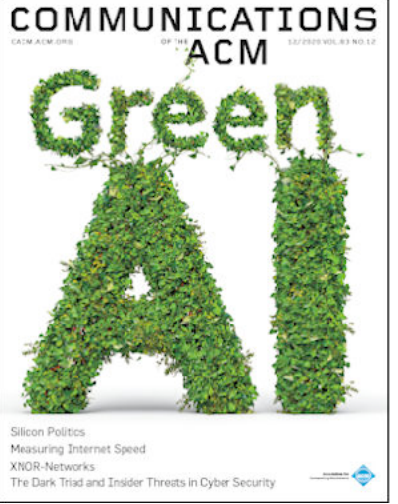
- What makes a good sparse structure?



Efficiency

Open Questions

- What makes a good sparse structure?
- Combining different methods



Think Green!

- **Red AI**

- Problems: inclusiveness, environment

- **Green AI**

- Enhance **reporting** of computational budgets

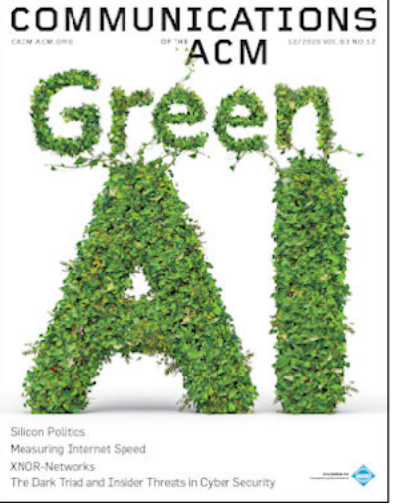


- Add a *price-tag* for scientific results

- Promote **efficiency** as a core evaluation for AI



- **In addition to** accuracy



Think Green!

- **Red AI**

- Problems: inclusiveness, environment

- **Green AI**

- Enhance **reporting** of computational budgets

- Add a *price-tag* for scientific results

- Promote **efficiency** as a core evaluation for AI

- **In addition to** accuracy

