

Generalization through Memorization: Nearest Neighbor Language Models

Authors: Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, Mike Lewis

Stanford University, Facebook AI Research

ICLR 2020

Presented by: Itay Levy

Talk outline

kNN-LM

- Motivation
- Constructing the datastore
- Inference

Experiment setting

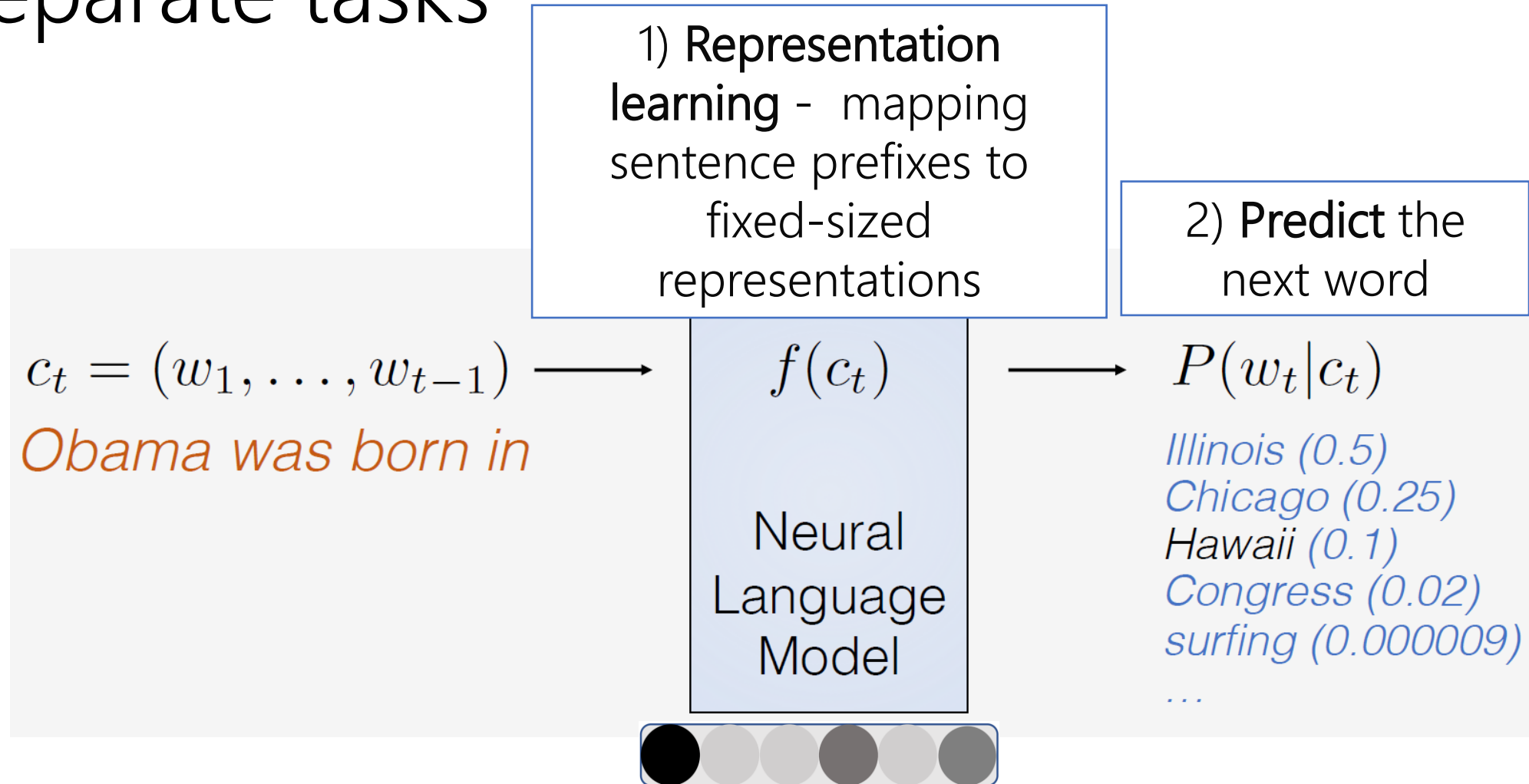
Key results

- Improved generalization w/o adding parameters
- Scaling up to larger training sets using the datastore
- Domain adaptation

Summary

- Maybe also some additional stuff

Motivation: Autoregressive LM as two separate tasks



Motivation: Hypotheses

1. Representation learning may be easier for LM than the prediction problem
2. Context representation function is a similarity function.
 - Contexts which are close in representation space are more likely to be followed by the same target word.
 - E.g. *"Dickens is the author of"* ~ *"Dickens wrote"*
 - *Seen before in the seminar* - in openQA we compare the query representation to the passage representation

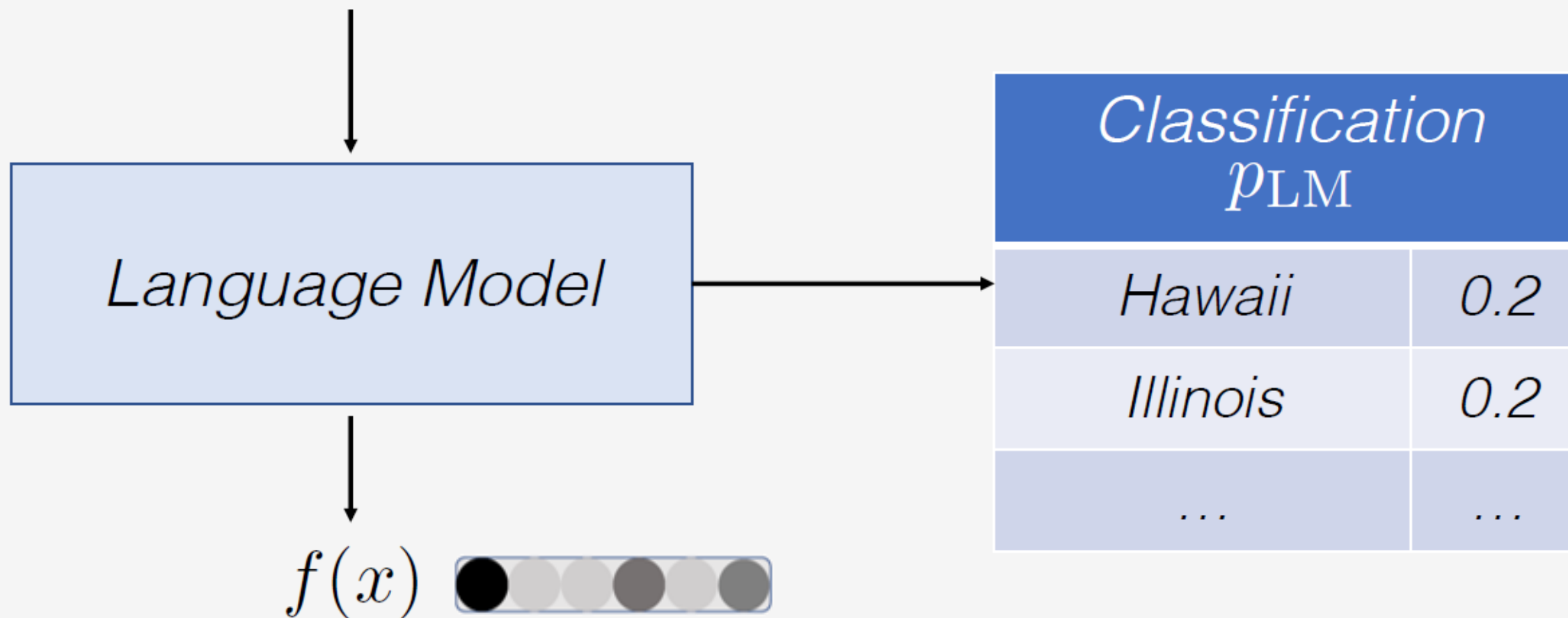
Motivation: finding close contexts in representation space

Test Context: Obama's birthplace is ???

<i>Previously Seen Contexts</i>	<i>Targets</i>
<i>Obama was senator for</i>	<i>Illinois</i>
<i>Barack is married to</i>	<i>Michelle</i>
<i>Obama was born in</i>	<i>Hawaii</i>
<i>...</i>	<i>...</i>
<i>Obama is a native of</i>	<i>Hawaii</i>


kNN-LM motivation

$x =$ *Obama's birthplace is _____*



kNN-LM motivation

$x =$ *Obama's birthplace is _____*

$q = f(x) =$ 




*Nearest Neighbors
Datastore*

Focus on neighbors for inference
efficiency (like in openQA)

kNN-LM motivation

$x =$ *Obama's birthplace is _____*

$q = f(x) =$ 



<u>Keys</u>	<u>Values</u>
<i>f(Obama was senator for)</i>	<i>Illinois</i>
<i>f(Obama was born in)</i>	<i>Hawaii</i>
...	...

Talk outline

kNN-LM



- Motivation
- Constructing the datastore
- Inference

Experiment
setting

Key results

- Improved generalization w/o adding parameters
- Scaling up to larger training sets using the datastore
- Domain adaptation

Summary

- Maybe also some additional stuff

Constructing the datastore





The text collection for the NN datastore can be the original LM training data or a different dataset

<i>Training Contexts</i> c_i	<i>Targets</i> v_i
<i>Obama was senator for</i>	<i>Illinois</i>
<i>Barack is married to</i>	<i>Michelle</i>
<i>Obama was born in</i>	<i>Hawaii</i>
<i>...</i>	<i>...</i>
<i>Obama is a native of</i>	<i>Hawaii</i>

Constructing the datastore

Create entry for each
training set token

Datastore

<i>Training Contexts</i> c_i	Keys	Values
	<i>Representations</i> $k_i = f(c_i)$	<i>Targets</i> v_i
<i>Obama was senator for</i>		<i>Illinois</i>
<i>Barack is married to</i>		<i>Michelle</i>
<i>Obama was born in</i>		<i>Hawaii</i>
...
<i>Obama is a native of</i>		<i>Hawaii</i>

Constructing the datastore - efficiency

Task	Time for Wikitext-103 [hours]	Hardware
[reference] A single epoch of LM training forward + backward pass	5	1 GPU
A single forward pass over the dataset to save keys/values	4	1 GPU(+SSD storage)
Creating the datastore using FAISS	2	1 CPU(+SSD storage)

Time complexity - roughly like 1 epoch of training

Highly parallelizable method

Datastore implementation details

How large is it?

- Key for every token
- 4000 bytes for each key
- **~400GB of storage**

FAISS -open source library for fast nearest neighbor retrieval in high dimensional spaces

- Reduces I/O
 - Clusters the keys
 - Search based on cluster centroids
- Reduces storage usage
 - Stores compressed versions of the vectors (64 bytes)

Talk outline

kNN-LM



- Motivation
- Constructing the datastore
- Inference

Experiment setting

Key results

- Improved generalization w/o adding parameters
- Scaling up to larger training sets using the datastore
- Domain adaptation





Summary

- Maybe also some additional stuff

Inference – neighbors retrieval

The k -nearest neighbors for $q = f(x)$

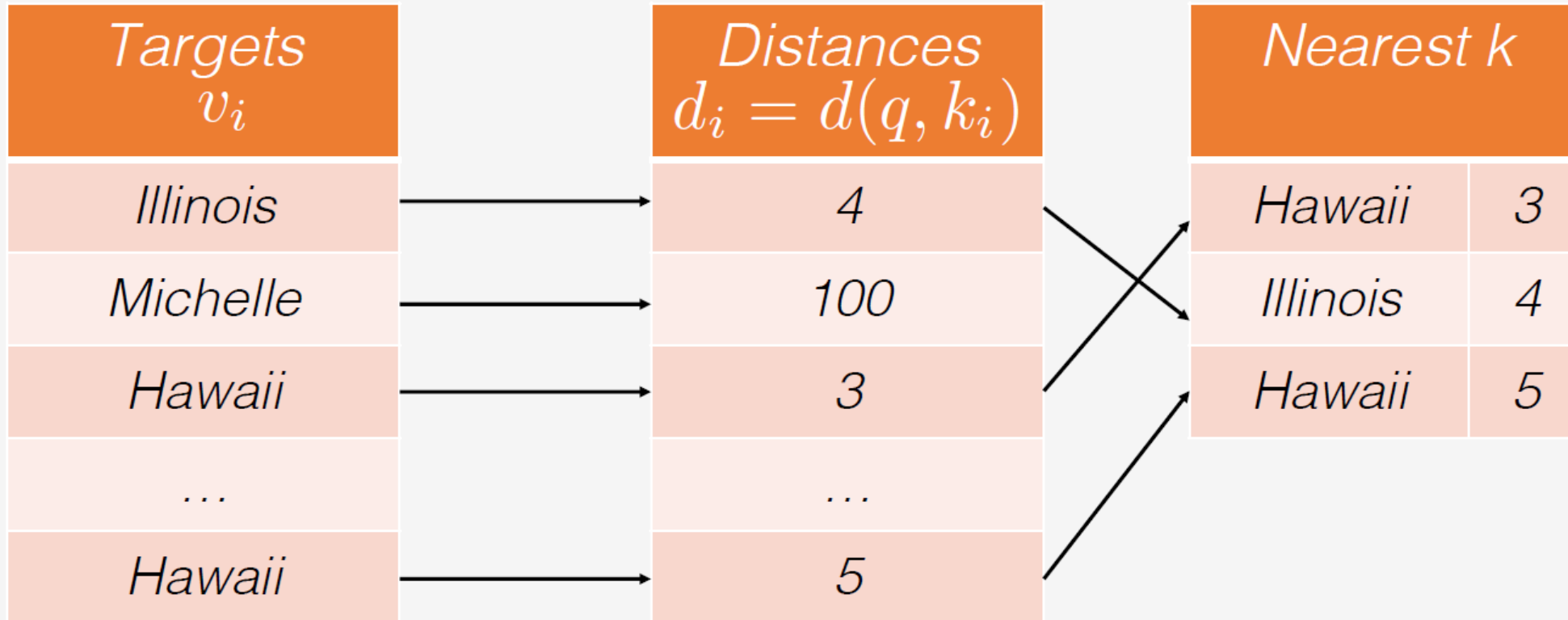


<i>Representations</i> $k_i = f(c_i)$	<i>Targets</i> v_i	<i>Distances</i> $d_i = d(q, k_i)$
	<i>Illinois</i>	4
	<i>Michelle</i>	100
	<i>Hawaii</i>	3
...
	<i>Hawaii</i>	5

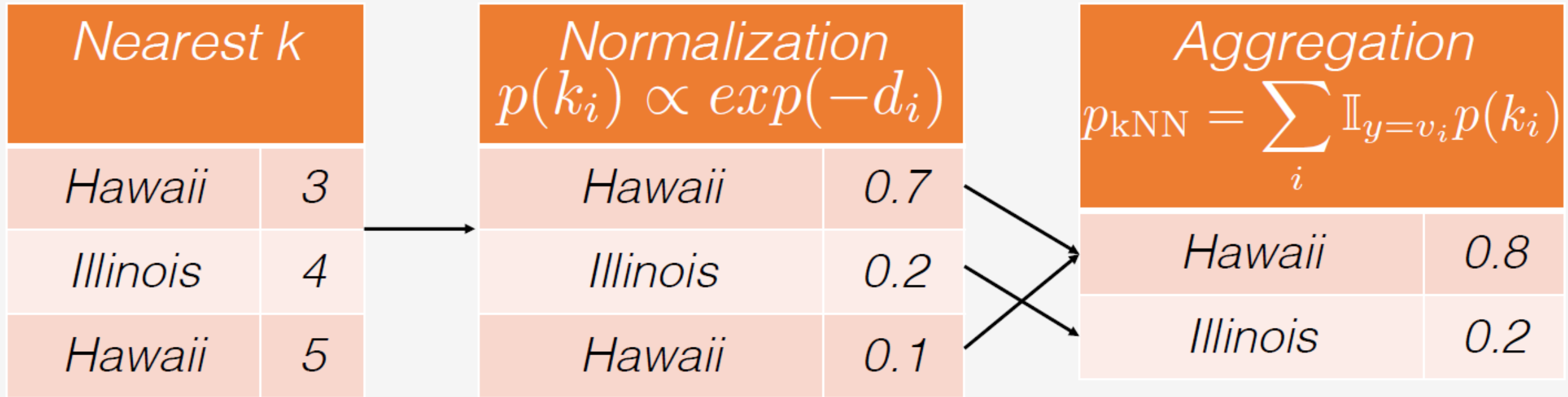
$K = 1024$

Inference – choosing nearest neighbors

The k -nearest neighbors for $q = f(x)$



Inference – kNN distribution



Softmax over negative distances

Inference – interpolation with the original model

$x =$ *Obama's birthplace is* _____

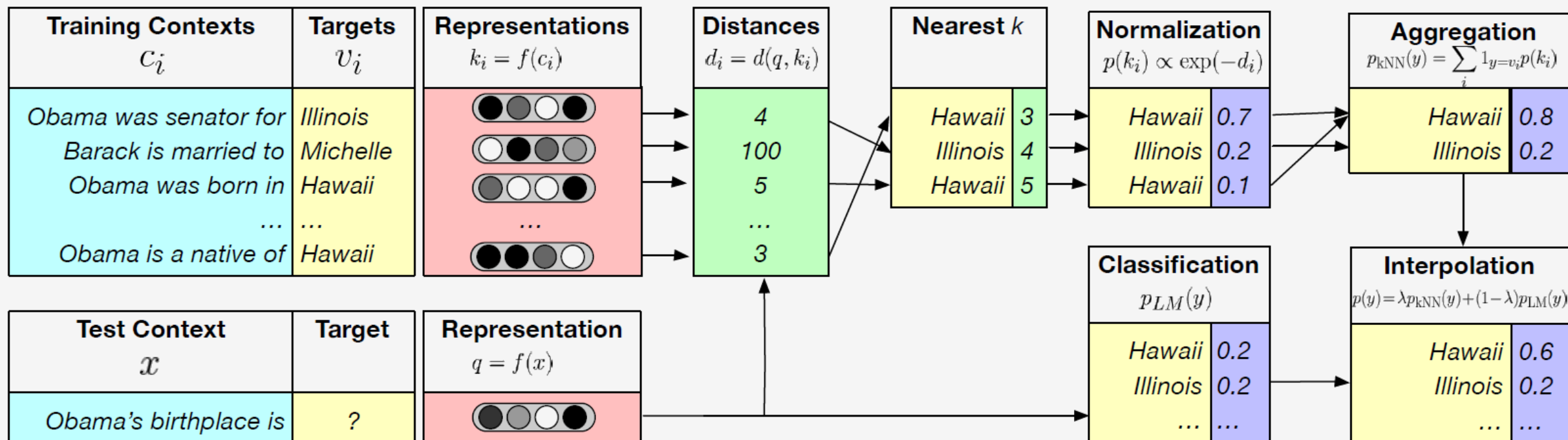
<i>Language Model</i>	
<i>Hawaii</i>	<i>0.2</i>
<i>Illinois</i>	<i>0.2</i>
<i>...</i>	<i>...</i>

<i>k-Nearest Neighbors</i>	
<i>Hawaii</i>	<i>0.8</i>
<i>Illinois</i>	<i>0.2</i>



<i>kNN-LM</i> $(1 - \lambda) p_{\text{LM}} + \lambda p_{\text{kNN}}$	
<i>Hawaii</i>	<i>0.6</i>
<i>Illinois</i>	<i>0.2</i>
<i>...</i>	<i>...</i>

Inference – complete illustration



Highly interpretable

Inference - efficiency

Language model	Decoding speed on a single GPU [tokens/sec]
Vanilla LM	500
kNN-LM	60

Talk outline

kNN-LM

- Motivation
- Constructing the datastore
- Inference



Experiment setting

Key results

- Improved generalization w/o adding parameters
- Scaling up to larger training sets using the datastore
- Domain adaptation

Summary

- Maybe also some additional stuff

Datasets

Dataset	Description	#Tokens
WIKITEXT-103	Standard benchmark for autoregressive language modelling Wikipedia articles	100M
BOOKS (Toronto Books Corpus)	10K books in 16 different genres	1B
WIKI-3B	English Wikipedia articles	3B
WIKI-100M	Random articles subset of WIKI-3B	100M

Model details

- Decoder-only Transformers
 - Architecture and optimization from Baevski & Auli (2019)
 - 16 layers
 - The keys in the datastore are the input to the final layer's feedforward network
 - 1024 dimensional hidden states – our keys

Evaluation

- Loss function during training
 - negative log-likelihood of the training corpus
- Evaluation metric
 - Perplexity (standard practice)

Metric	Definition	Formula
Likelihood	$P(D)$	$\prod_{y_i \in D} P(y_i Y_{<i})$
Log-likelihood	$\log P(D)$	$\sum_{y_i \in D} \log P(y_i Y_{<i})$
Cross Entropy	$-\frac{1}{ D } \log P(D)$	$-\frac{1}{ D } \sum_{y_i \in D} \log P(y_i Y_{<i})$
Perplexity	$(P(D))^{-\frac{1}{ D }}$	$\exp\left(-\frac{1}{ D } \sum_{y_i \in D} \log P(y_i Y_{<i})\right)$

Credit: Advanced Methods in NLP course by Omer Levy

Qualitative Analysis

Memorizing rare information

The model is specifically helpful in predicting rare patterns

- factual knowledge
- Names
- near-duplicate sentences from the training set

Test Context ($p_{\text{kNN}} = 0.998, p_{\text{LM}} = 0.124$)	Test Target	
<i>it was organised by <u>New Zealand</u> international <u>player</u> <u>Joseph Warbrick</u>, promoted by civil servant <u>Thomas Eyton</u>, and <u>managed</u> by <u>James Scott</u>, a publican. The <u>Natives</u> were the first New Zealand team to perform a haka, and also the first to wear all black. They played 107 rugby matches during the tour, as well as a small number of Victorian Rules football and association football matches in Australia. Having made <u>a significant impact on the...</u></i>	development	
Training Set Context	Training Set Target	Context Probability
<i>As the captain and instigator of the 1888-89 <u>Natives</u> – the first <u>New Zealand</u> team to tour the British Isles – <u>Warbrick</u> had a lasting <u>impact on the...</u></i>	development	0.998
<i>promoted to a new first grade competition which started in 1900. <u>Glebe</u> immediately made a big <u>impact on the...</u></i>	district	0.00012
<i>centuries, few were as large as other <u>players</u> <u>managed</u>. However, others contend that his <u>impact on the...</u></i>	game	0.000034
<i>Nearly every game in the main series has either an anime or manga adaptation, or both. The series has had <u>a significant impact on the...</u></i>	development	0.00000092

Figure 6: Example where the k NN model has much higher confidence in the correct target than the LM. Although there are other training set examples with similar local n -gram matches, the nearest neighbour search is highly confident of specific and very relevant context.

Talk outline

kNN-LM

Experiment
setting

Key results

Summary



sentiment
analysis

vibe
check



- Improved generalization w/o adding parameters
- Scaling up to larger training sets using the datastore
- Domain adaptation

- Maybe also some additional stuff

Using the training data as the datastore

<i>Model</i>	<i>Perplexity</i>
<i>Previous Best (Luo et al., 2019)</i>	<i>17.40</i>
<i>Base LM</i>	<i>18.65</i>
<i>kNN-LM</i>	<i>16.12</i>



Very significant result!

Key result #1

Improved generalization w/o adding parameters

- Compatible with any autoregressive model
- No added parameters*
- No additional training
 - Takes advantage of effective similarity functions learned by LM

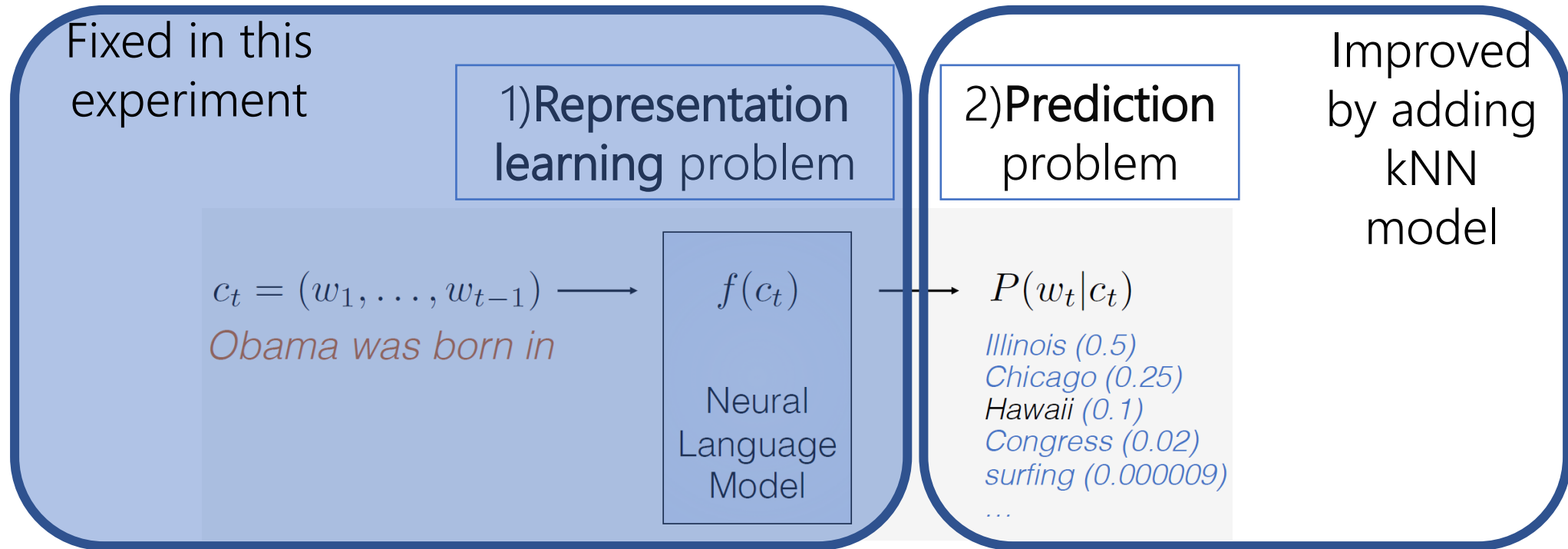
Is it just ensembling with an overfitted model?

Memorizing Transformer – Transformer LM that perfectly memorized the training data

Interpolating original LM with:	Validation perplexity improvement ↑
kNN-LM	1.9
Memorizing Transformer	0.1

Implicit memorization is less effective at generalization than kNN-LM

Back to the hypotheses



1. Prediction problem is harder when done implicitly using LM parameters
2. Contexts which are close in representation space are more likely to be followed by the same target word

Talk outline

kNN-LM

- Motivation
- Constructing the datastore
- Inference

Experiment setting

Key results

- Improved generalization w/o adding parameters
- Scaling up to larger training sets using the datastore
- Domain adaptation

Summary

- Maybe also some additional stuff

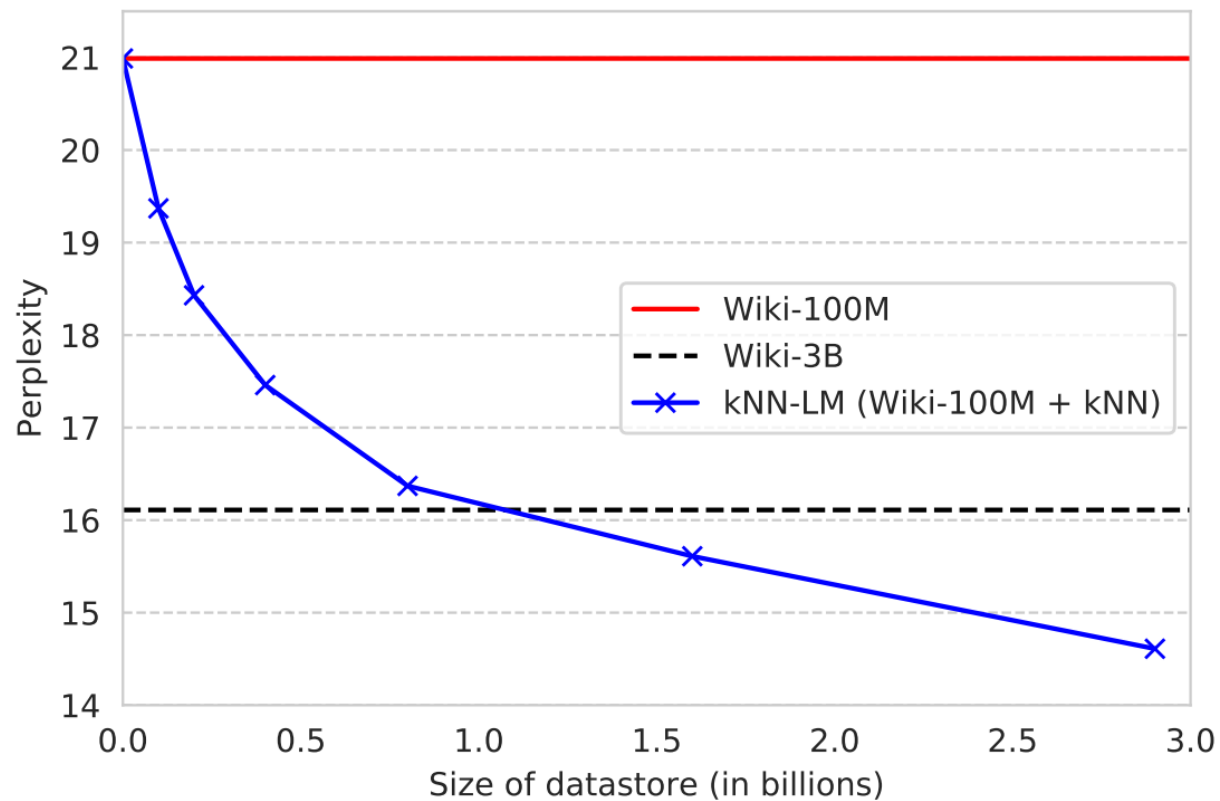
What should I do with extra data?

<i>LM Training Data</i>	<i>Datastore</i>	<i>Perplexity</i>
<i>Wiki-3B</i>	-	<i>15.17</i>
<i>Wiki-100M</i>	-	<i>19.59</i>
<i>Wiki-100M</i>	<i>Wiki-3B</i>	<i>13.73</i>

Test on Wiki-3B

Retrieving nearest neighbors from the corpus outperforms training on it!

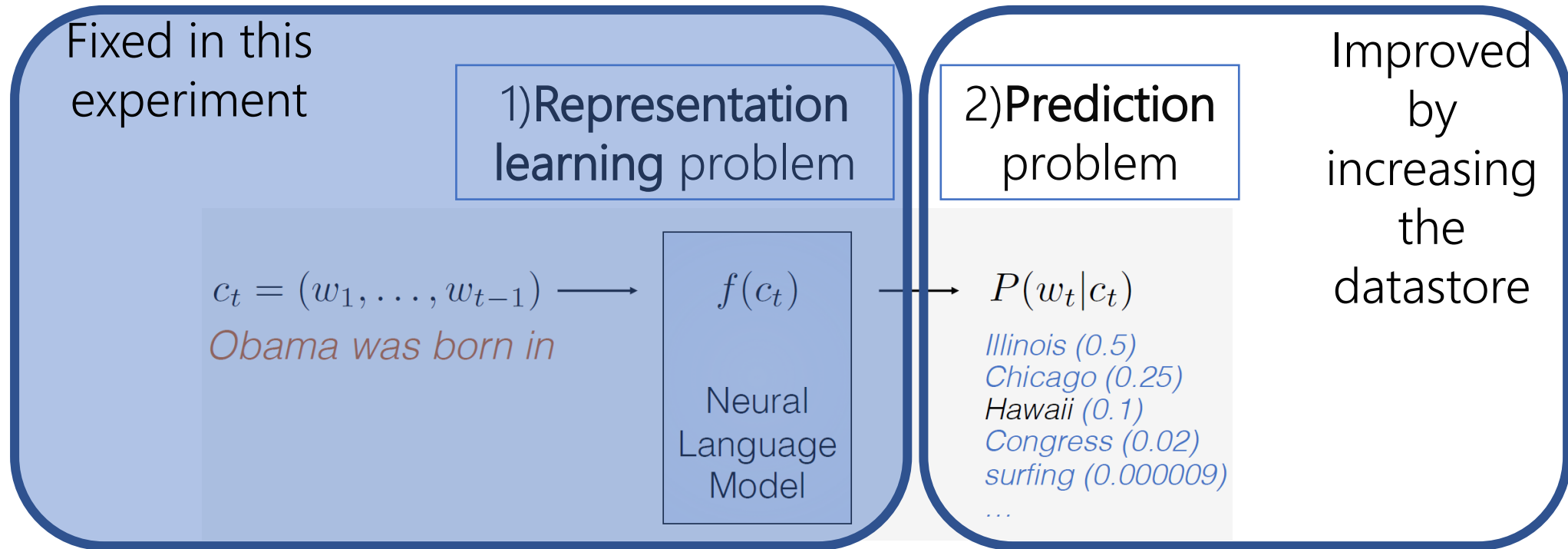
What should I do with extra data?



(a) Effect of datastore size on perplexities.

kNN-LM trained on 100M tokens with a datastore of 1.6B tokens already outperforms the LM trained on all 3B tokens

Back to the hypotheses



Prediction problem is harder when done implicitly using LM parameters

Possible implication?

Massive LM (like GPT-3) are not that better at finding better representations

They are better at prediction by implicitly memorizing more data in parameters

Credit: Advanced Methods in NLP course by Omer Levy

Key result #2

Scaling up to larger training sets using the datastore

- Retrieving nearest neighbors from the corpus outperforms training on it
- New path for efficiently using large datasets in LMs
 - No additional cost of training
 - Just increasing the datastore
- Problem
 - Larger datastore -> slower inference

Talk outline

kNN-LM

- Motivation
- Constructing the datastore
- Inference

Experiment setting

Key results



- Improved generalization w/o adding parameters
- Scaling up to larger training sets using the datastore
- Domain adaptation

Summary

- Maybe also some additional stuff

Domain adaptation

- Varying the NN datastore, again without further training
 - Adding out-of-domain data to the datastore makes a single LM useful across multiple domains

<i>LM Training Data</i>	<i>Datastore</i>	<i>Perplexity on Books</i>
<i>Books</i>	-	<i>11.89</i>
<i>Wiki-3B</i>	-	<i>34.84</i>
<i>Wiki-3B</i>	<i>Books</i>	<i>20.47</i>

Domain adaptation requires more weight on kNN component than in-domain

Key result #3

Domain adaptation

- A single LM can adapt to multiple domains without the in-domain training
- We can domain-specific data to the datastore

Talk outline

kNN-LM

- Motivation
- Constructing the datastore
- Inference

Experiment setting

Key results

- Improved generalization w/o adding parameters
- Scaling up to larger training sets using the datastore
- Domain adaptation



Summary

- Maybe also some additional stuff

Summary

- Explicitly memorizing the training data helps generalization
 - Suggests that prediction problem is harder than representation learning
 - Takes advantage of effective similarity functions learned by LM
- Enables:
 - Compatibility to any autoregressive LM
 - Smaller models trained on smaller datasets
 - Quicker training
 - Adaptability to other domains

Future work

- Explicitly training similarity functions.
- Reducing the size of the datastore.
- Keep training the decoding layers while freezing the encoding part that kNN-LM uses.
- Extend the proposed model to other language tasks.
- Test on various language models to verify the generalization ability across different models.

Result #1 on multiple domains

Control for the possibility that encyclopedic Wikipedia text is somehow uniquely good for caching.

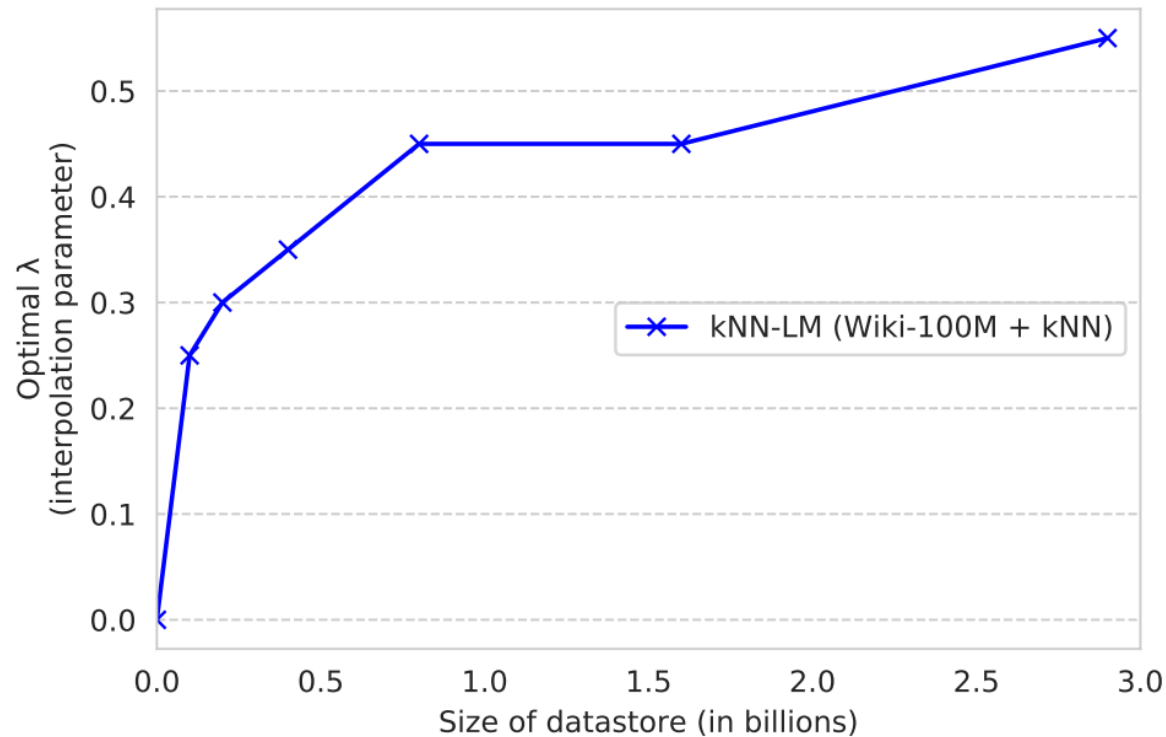
Model	Perplexity (\downarrow)		# Trainable Params
	Dev	Test	
Base LM (Baevski & Auli, 2019)	14.75	11.89	247M
+ k NN-LM	14.20	10.89	247M

Table 2: Performance on BOOKS, showing that k NN-LM works well in multiple domains.

Finding optimal λ

Datastore size

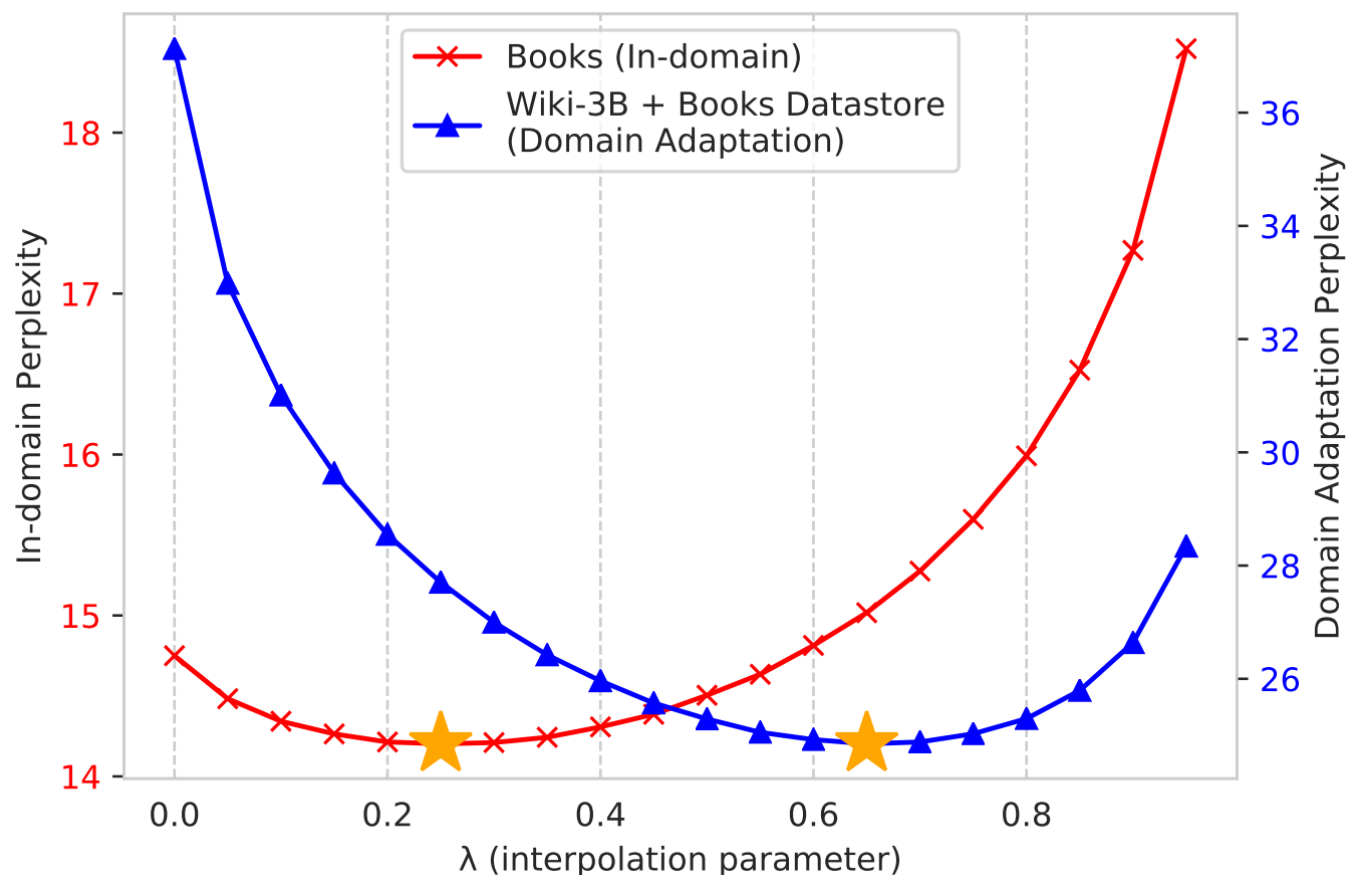
The optimal value of λ increases with the size of the datastore.



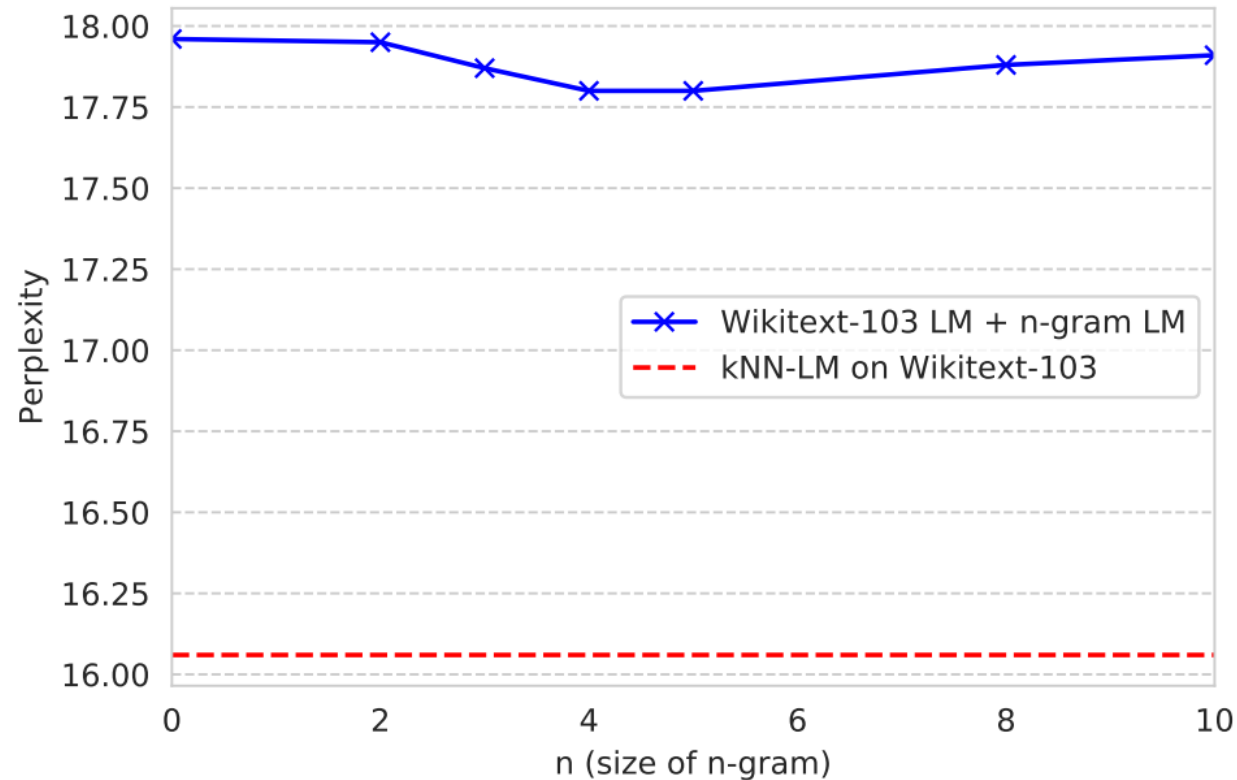
(b) Tuned values of λ for different datastore sizes.

Finding optimal λ in domain vs. out of domain

More weight on kNN LM improves domain adaptation.



kNN(neural representation) vs n-gram (simple representation)



Highlights the need to use the learned representation function $f(\cdot)$ to measure similarity between more varied contexts