

# Generalization through Memorization: Nearest Neighbor Language Models

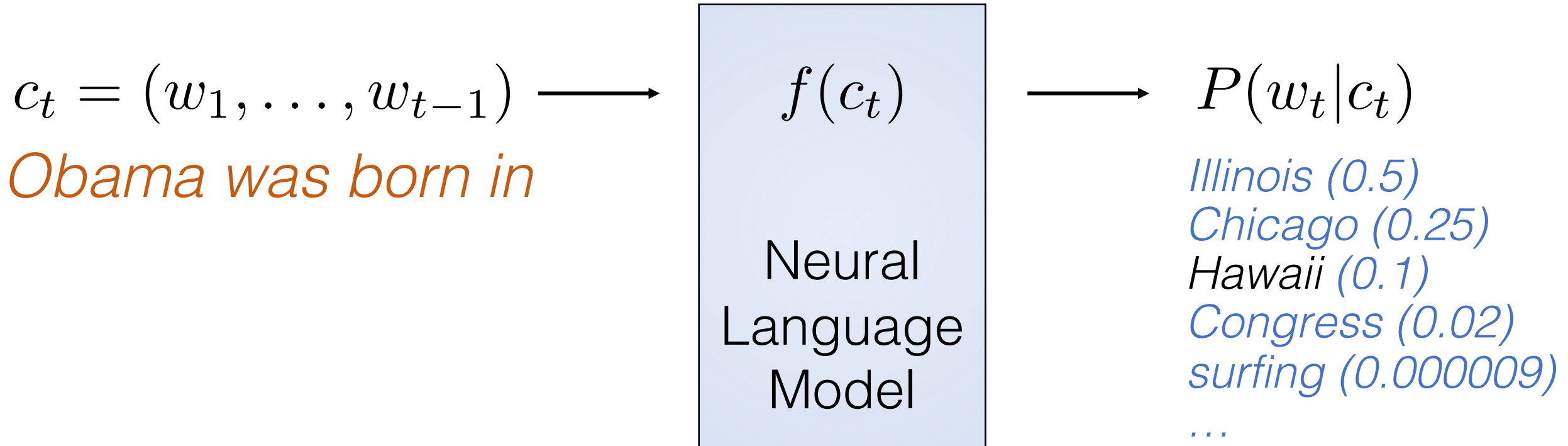
Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, Mike Lewis  
Stanford University, Facebook AI Research



**facebook** Artificial Intelligence

# Neural Autoregressive Language Models

*Given prior context, estimate the probability for the target token*



# Language Models

*Lots of text is very easily available, so we train models on large amounts of data.*

*But improving LM performance or scaling to larger datasets, by training bigger and bigger models with billions of parameters, requires massive amounts of GPU compute.. 😞*

*Instead, can explicitly **memorizing** data make LMs **generalize** better without the added cost of training?*

# Nearest Neighbor Language Models



# Key Results



*Explicitly memorizing* the training data helps generalization.

LMs can *scale* to larger text collections without the added cost of training.

A single LM can *adapt* to multiple domains without any in-domain training.

# Nearest Neighbor Language Models (kNN-LM)

# kNN-LM: Intuition

*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>



Given a new test context...

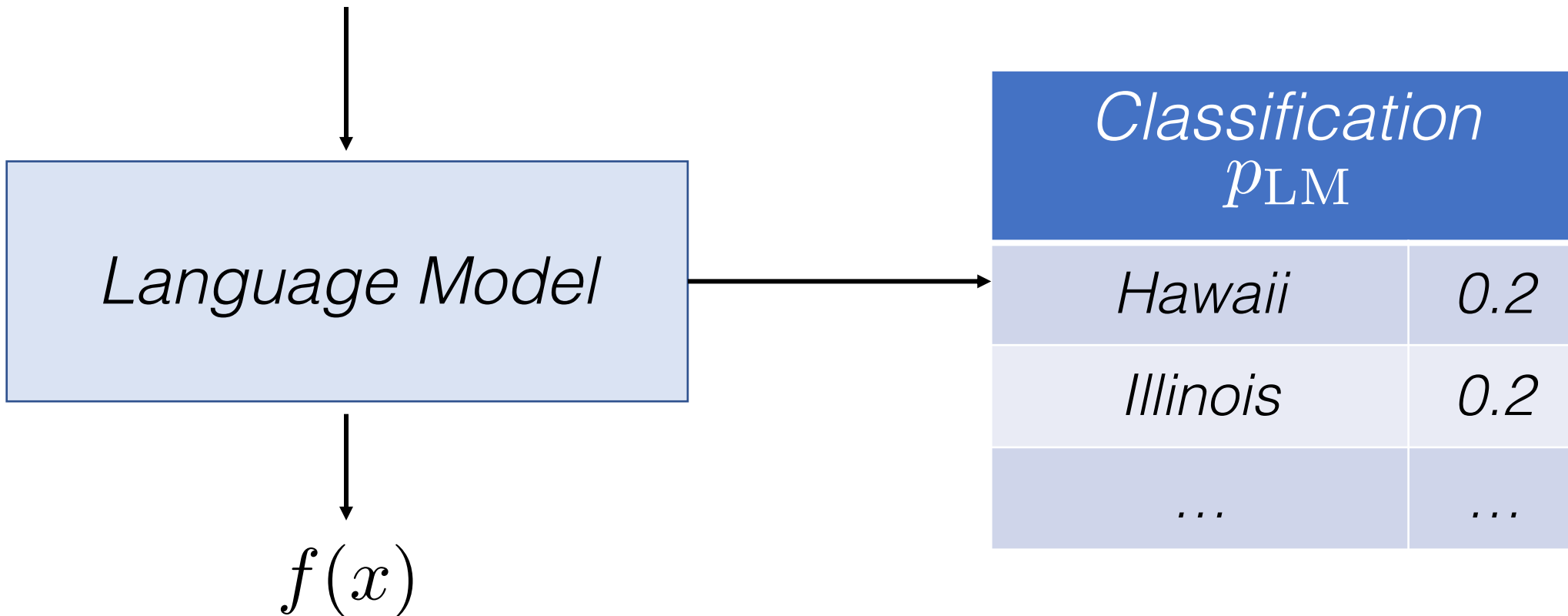
$x =$  *Obama's birthplace is* \_\_\_\_\_



*Language Model*

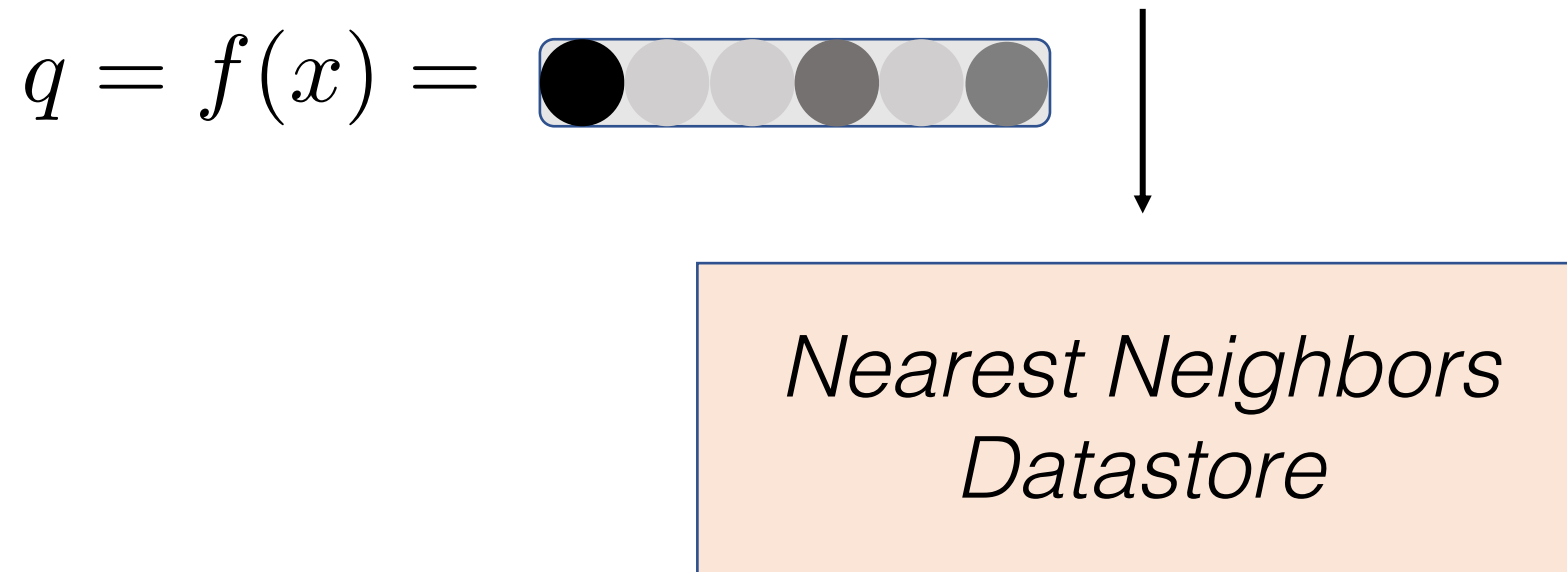
Given a new test context...

$x =$  *Obama's birthplace is \_\_\_\_\_*




Given a new test context...

$x =$  *Obama's birthplace is* \_\_\_\_\_



Given a new test context...

$x =$  *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>
...	...

# Constructing the datastore

# Constructing the datastore

<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

<i>Training Contexts</i> $c_i$	<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>

Back to inference!

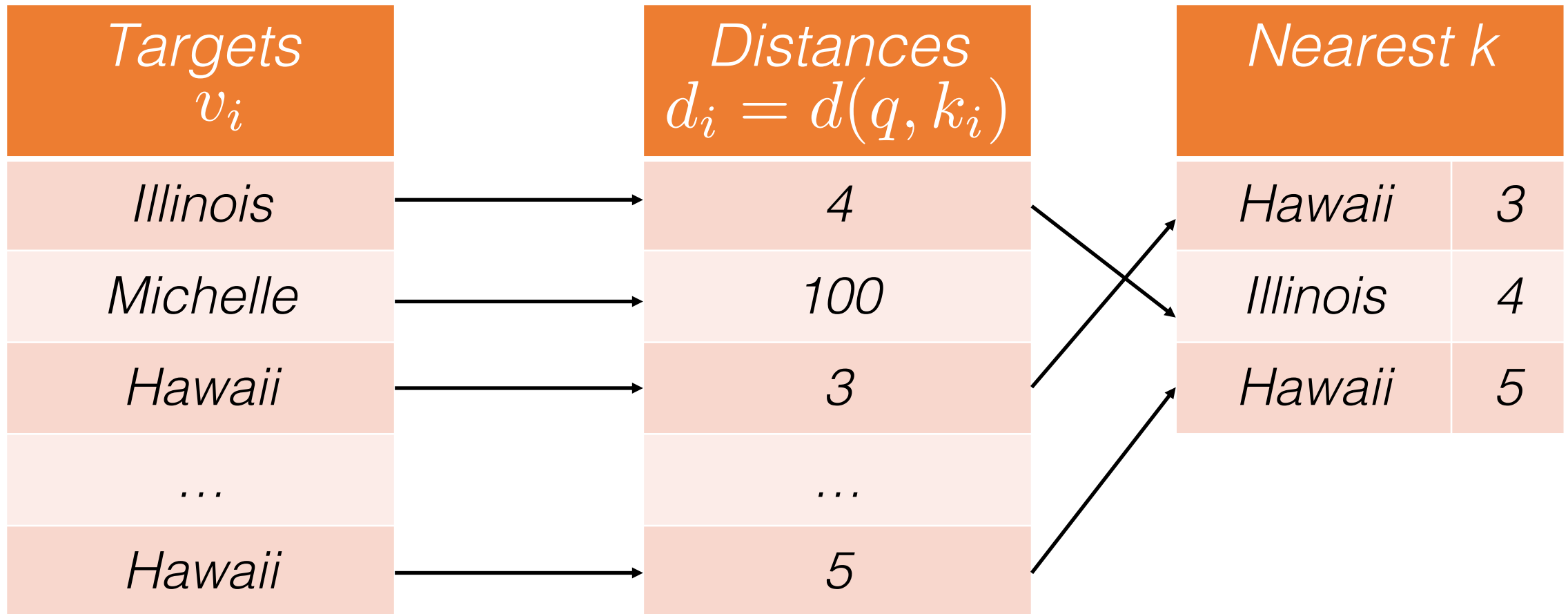


# The k-nearest neighbors for $q = f(x)$



<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

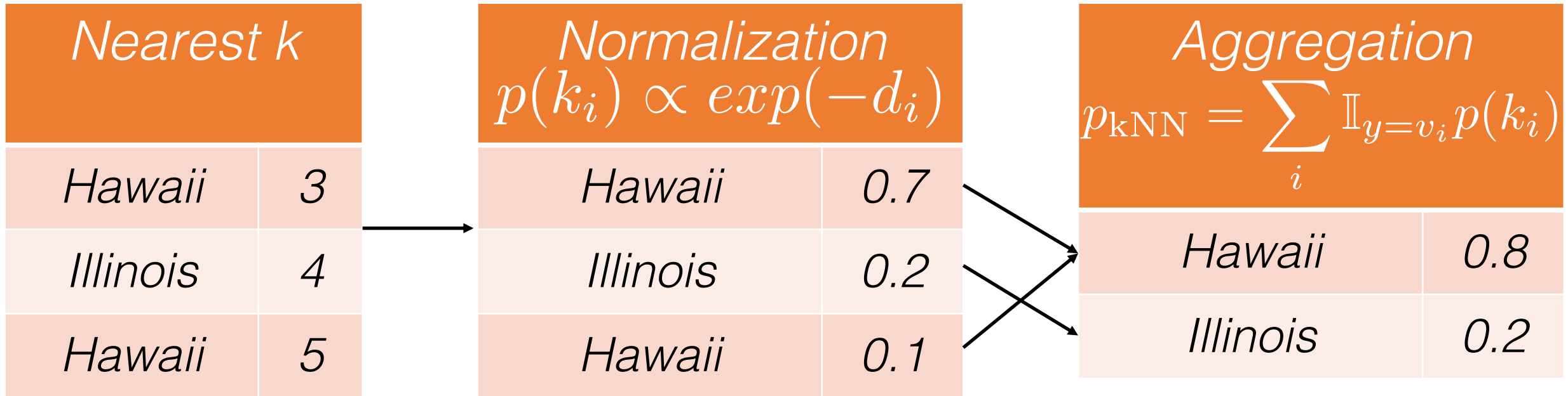
# The k-nearest neighbors for $q = f(x)$



# The kNN distribution

<i>Nearest k</i>			<i>Normalization</i> $p(k_i) \propto \exp(-d_i)$	
<i>Hawaii</i>	3	→	<i>Hawaii</i>	0.7
<i>Illinois</i>	4		<i>Illinois</i>	0.2
<i>Hawaii</i>	5		<i>Hawaii</i>	0.1

# The kNN distribution



Given a new test context...

$x =$  *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>



# Experiments

*Our Base LM is the **Transformer LM** from Baevski and Auli (2019).*

# Key Results



*Explicitly memorizing* the training data helps generalization.

*LMs can scale to larger text collections without the added cost of training, by simply adding the data to the datastore.*

*A single LM can adapt to multiple domains without the in-domain training, by adding domain-specific data to the datastore.*

# Memorizing with Wikitext-103

*Standard LM benchmark, 103 million tokens*

<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>



# Memorizing with Wikitext-103

*Datastore contains 103M examples,  $\lambda = 0.25$*

<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>



# Memorizing with Wikitext-103

*Datastore contains 103M examples,  $\lambda = 0.25$*

<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>
<i>kNN-LM + Cont. Cache*</i>	<i>15.79</i>



# Key Results



*Explicitly memorizing the training data helps generalization.*

*LMs can **scale** to larger text collections without the added cost of training, by simply adding the data to the datastore.*

*A single LM can adapt to multiple domains without the in-domain training, by adding domain-specific data to the datastore.*

# Scaling up from Wiki-100M to Wiki-3B

<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>

# Scaling up from Wiki-100M to Wiki-3B

<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>

# Scaling up from Wiki-100M to Wiki-3B

<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>

*Retrieving nearest neighbors from the corpus outperforms training on it!*

# Key Results



*Explicitly memorizing the training data helps generalization.*

*LMs can scale to larger text collections without the added cost of training, by simply adding the data to the datastore.*

*A single LM can **adapt** to multiple domains without the in-domain training, by adding domain-specific data to the datastore.*

# Domain Adaptation from Wiki to Books

<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>



# Domain Adaptation from Wiki to Books

<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 from Wiki to Books

<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>

*A single LM can be useful in multiple domains by simply adding a domain-specific datastore!*

# kNN-LM



*Explicitly memorizing* the training data helps generalization.

LMs can *scale* to larger text collections without the added cost of training, by simply adding the data to the datastore.

A single LM can *adapt* to multiple domains without the in-domain training, by adding domain-specific data to the datastore.

# Thanks!

*Explicitly memorizing* the training data helps generalization.

*LMs can scale* to larger text collections without the added cost of training, by simply adding the data to the datastore.

*A single LM can adapt* to multiple domains without the in-domain training, by adding domain-specific data to the datastore.



"To make a long story short, what it all boils down to in the final analysis is that what you should take away from this is..."

Paper:

<https://arxiv.org/pdf/1911.00172.pdf>

Code:

<https://github.com/urvashik/knnlm>