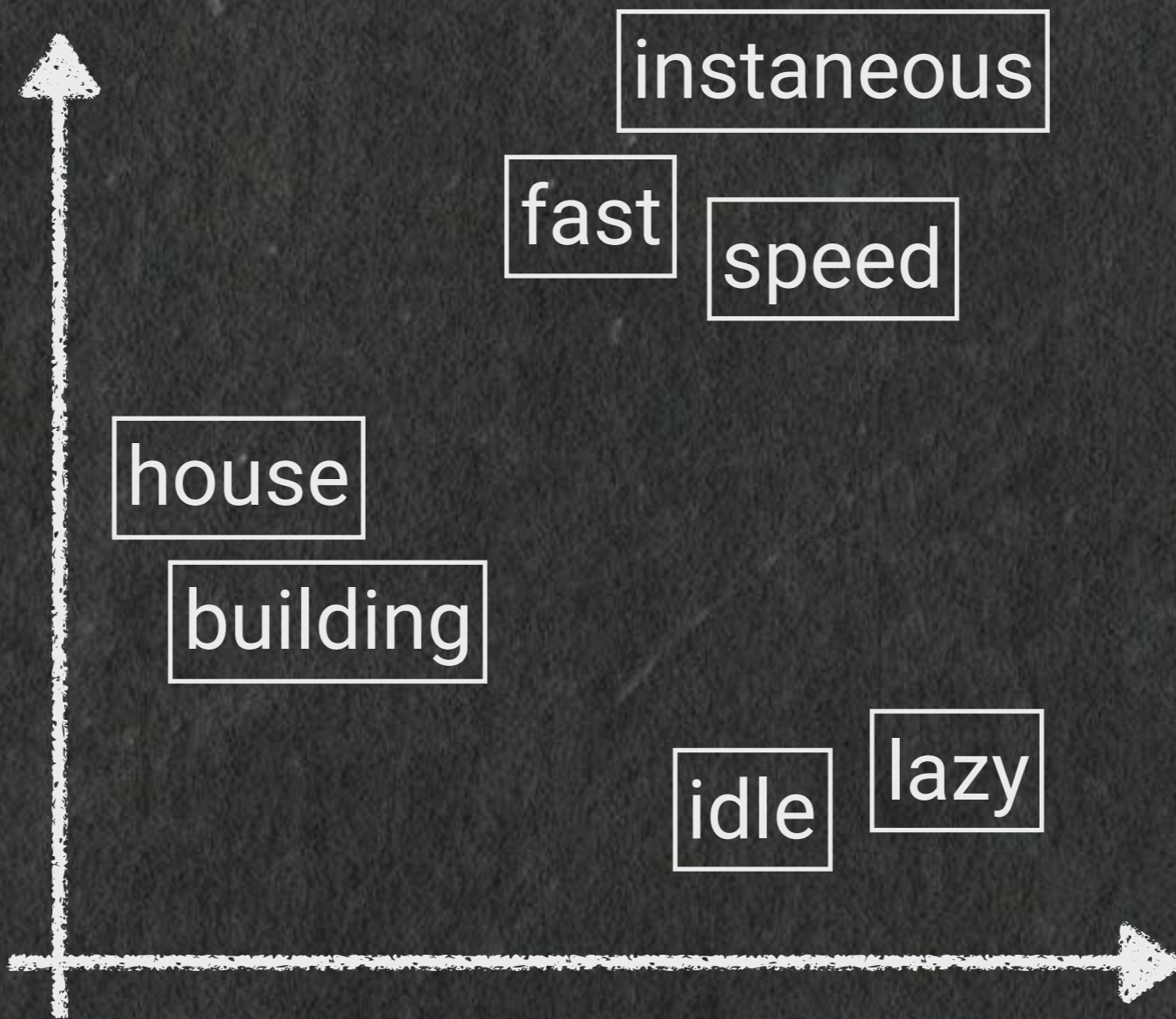


Feeding Word2vec
with tens of billions of items,
what could possibly go wrong?

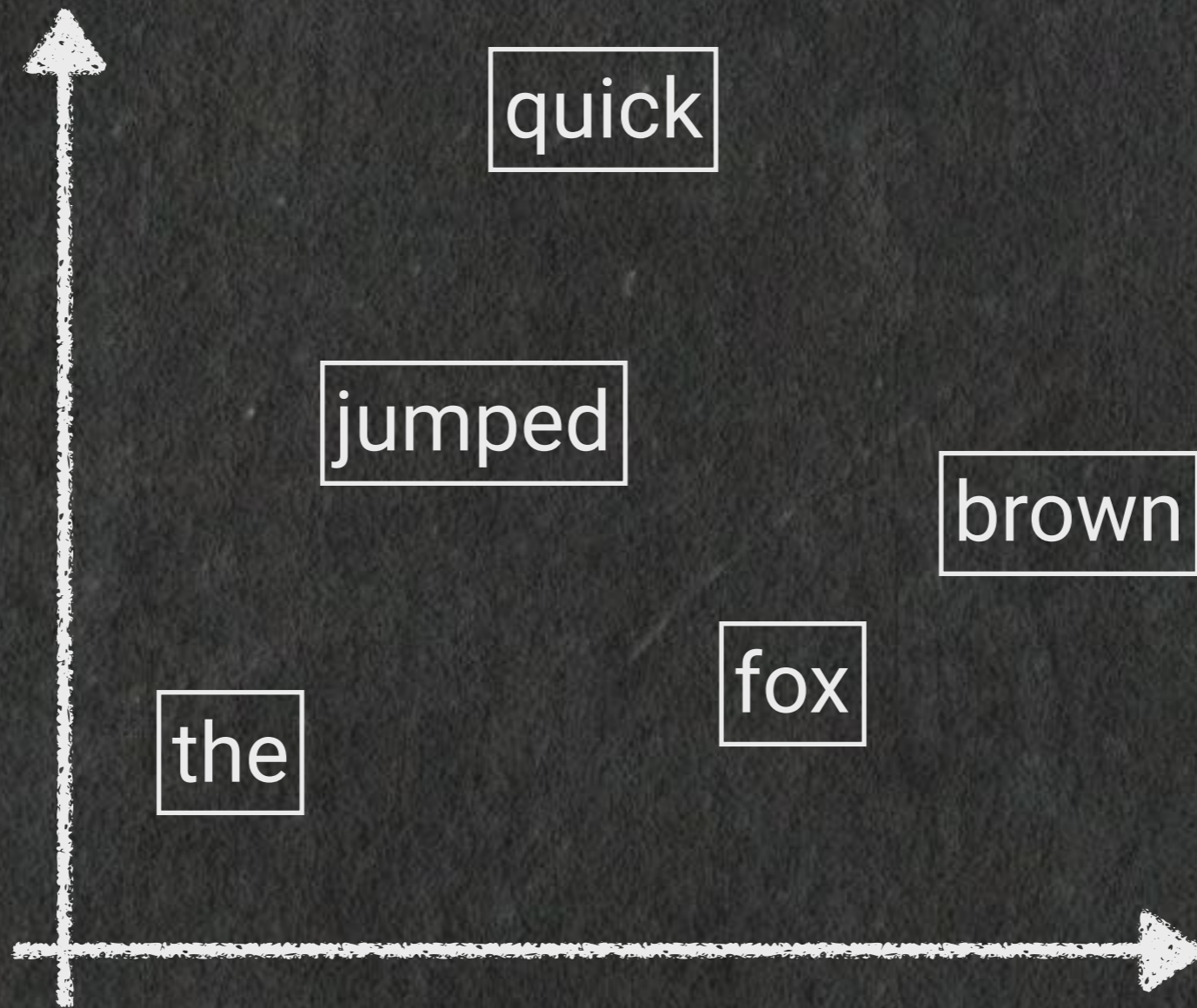
Simon Dollé, 12/06/2017





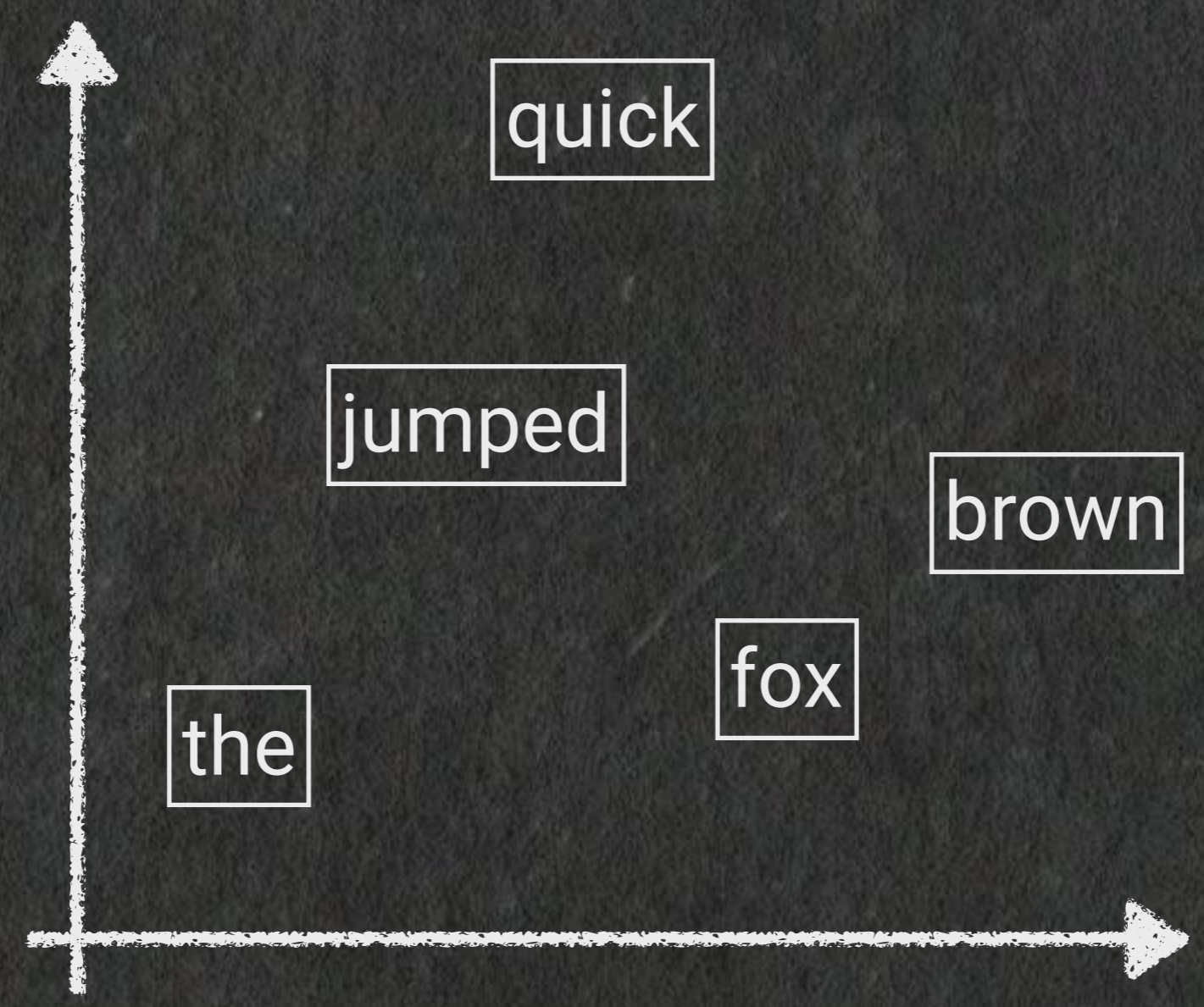
the quick brown fox jumped over the lazy dog

the quick brown fox jumped over the lazy dog



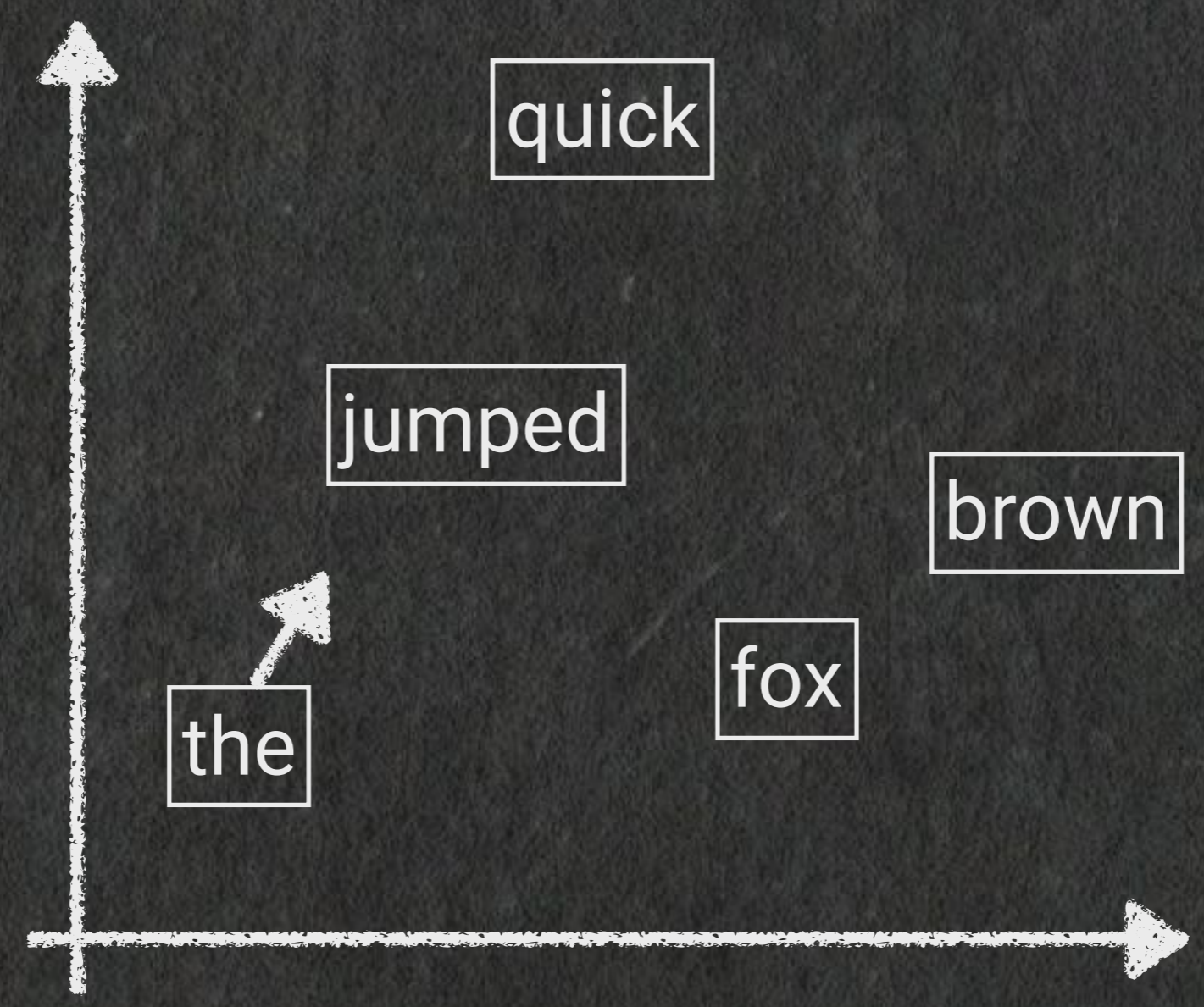
~~the~~ quick ~~brown~~ fox jumped over the lazy dog

? ?



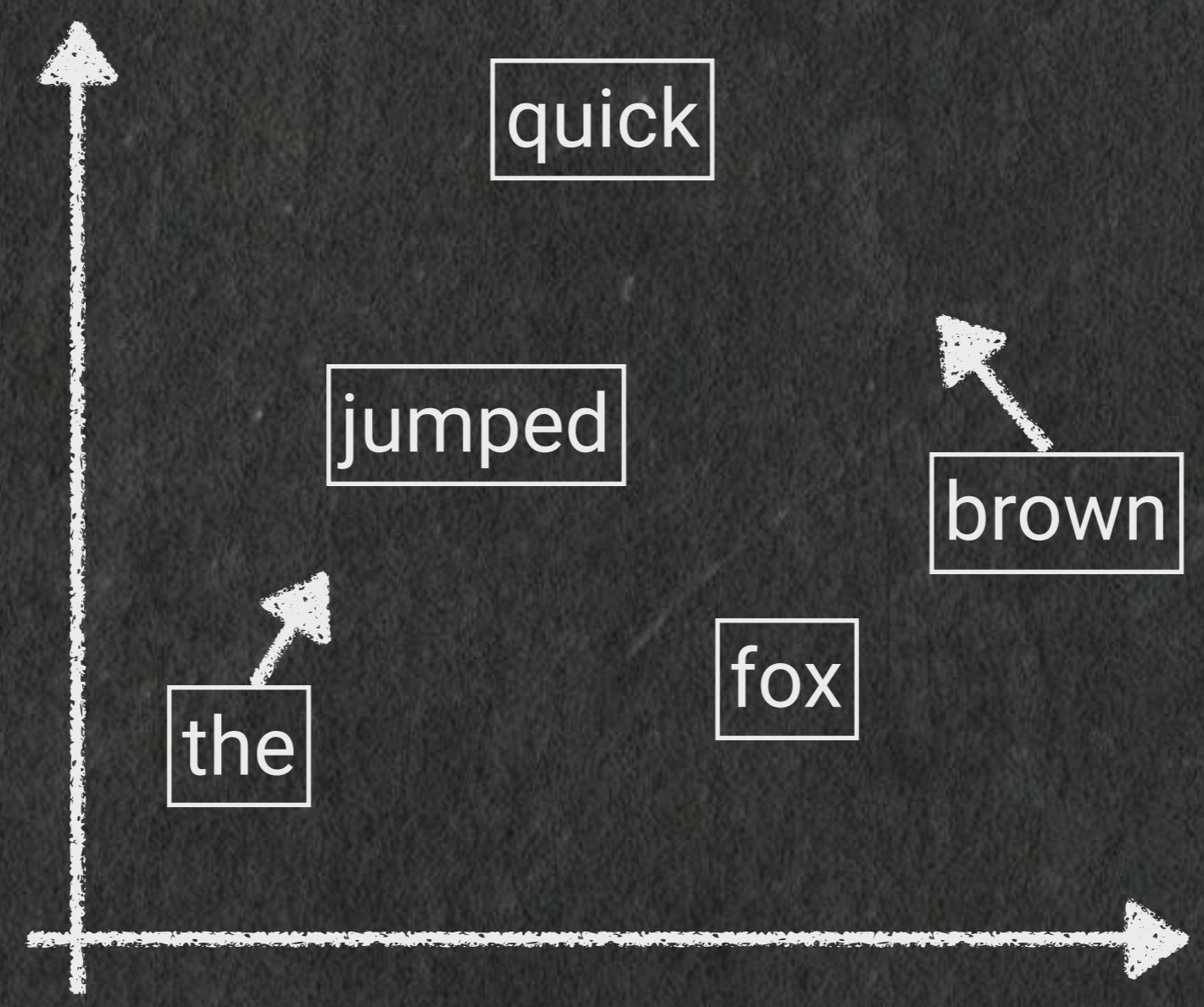
~~the~~ quick ~~brown~~ fox jumped over the lazy dog

? ?



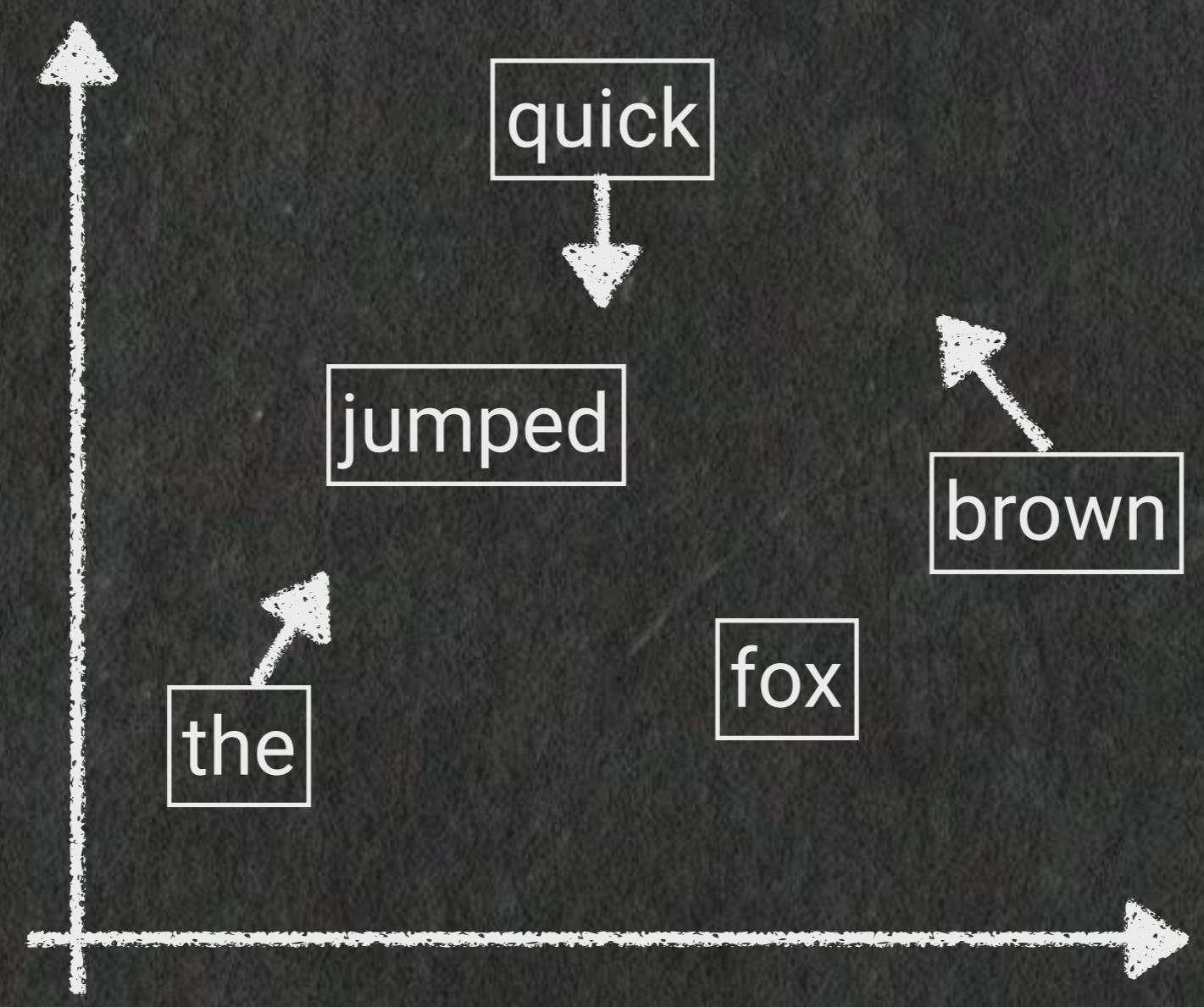
~~the~~ quick ~~brown~~ fox jumped over the lazy dog

? ?



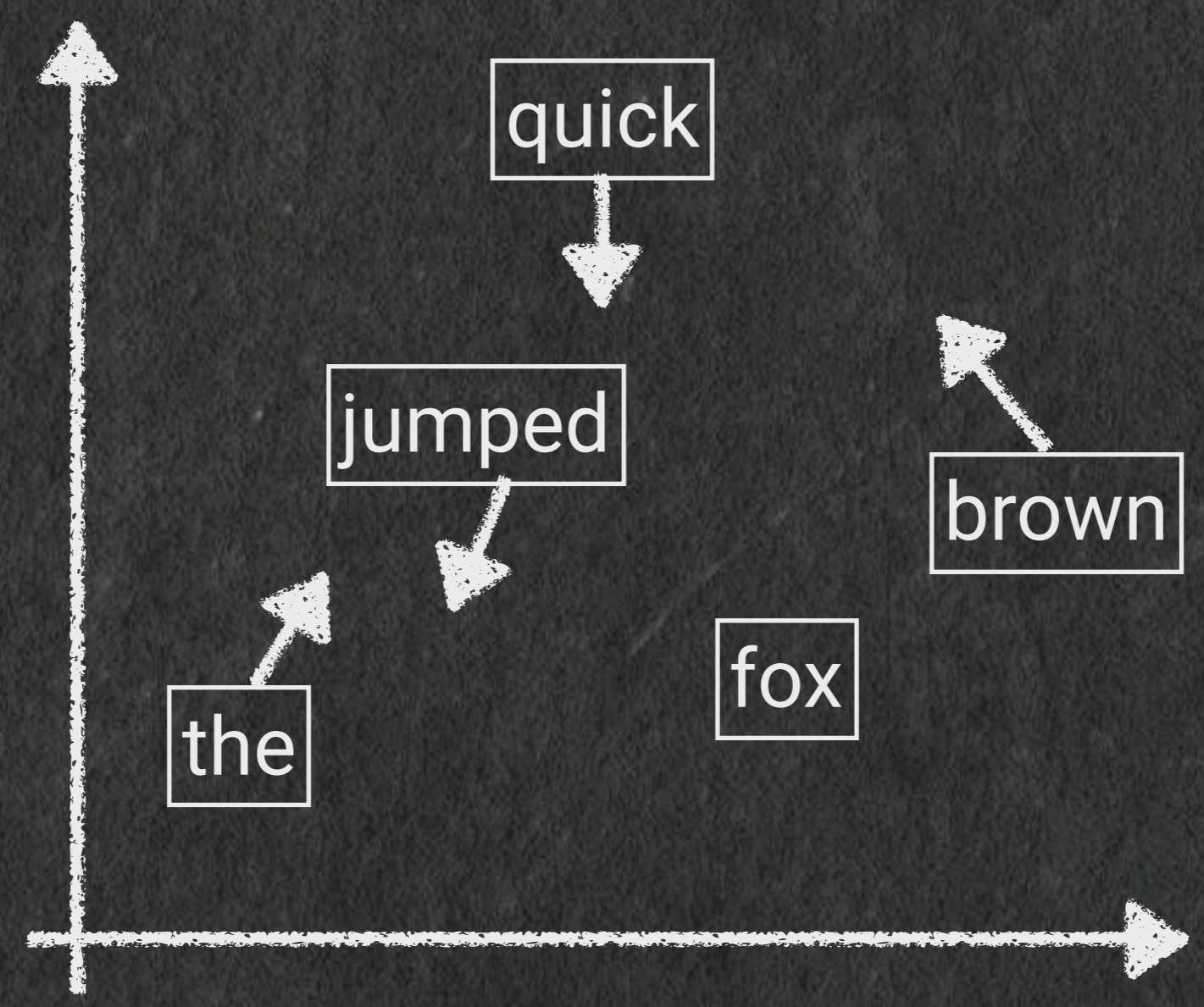
~~the~~ quick ~~brown~~ fox jumped over the lazy dog

? ?

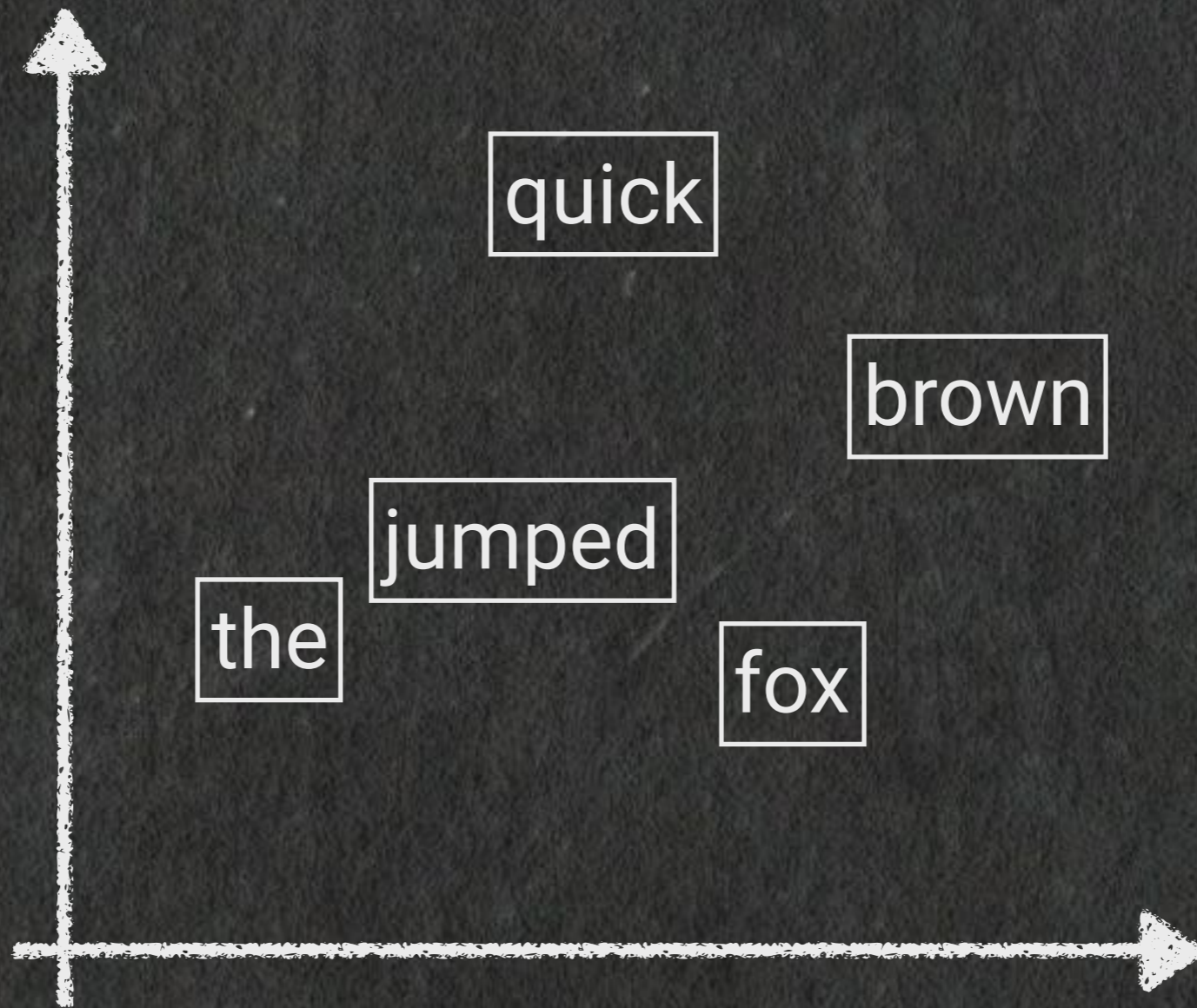


~~the~~ quick ~~brown~~ fox jumped over the lazy dog

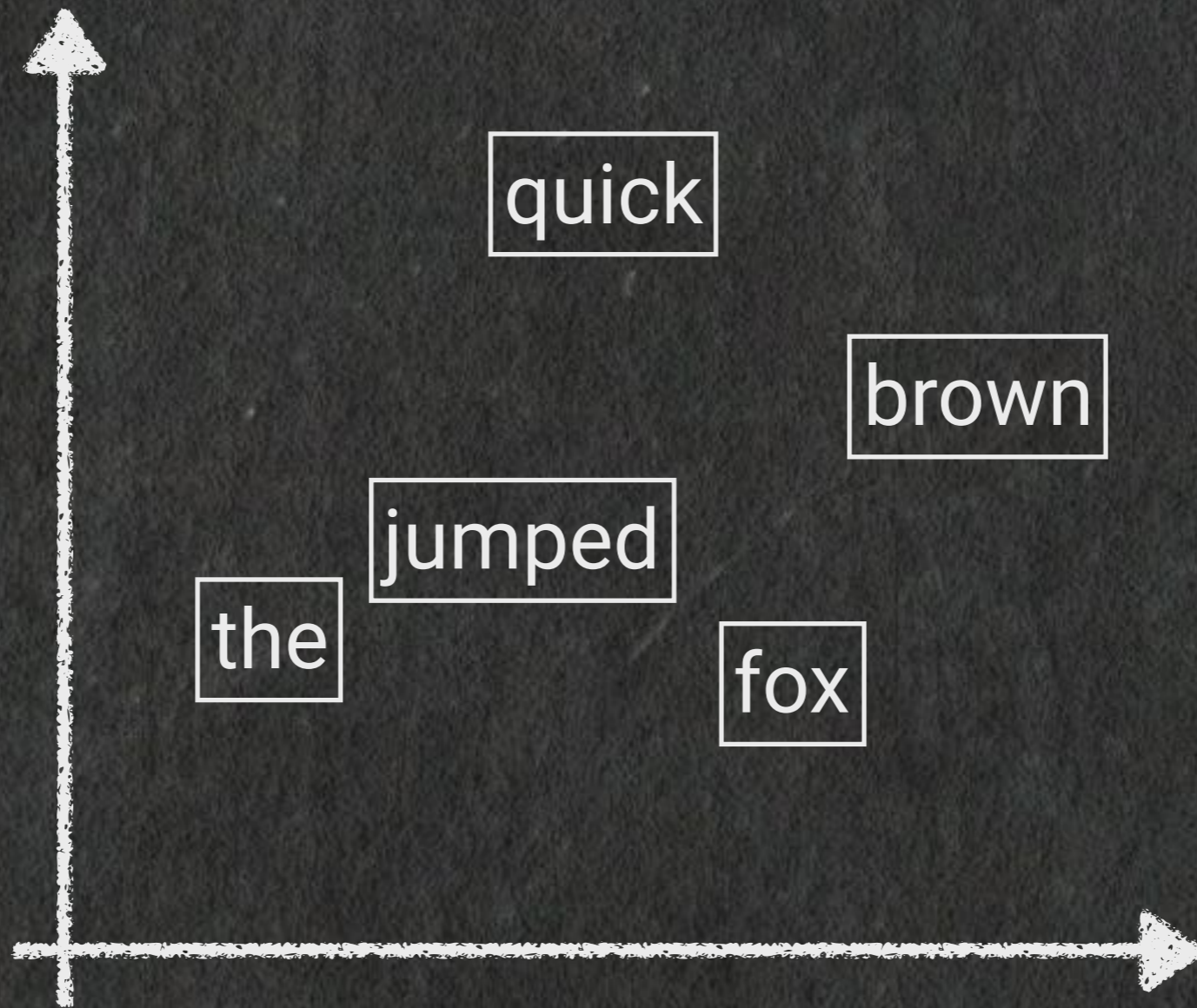
? ?



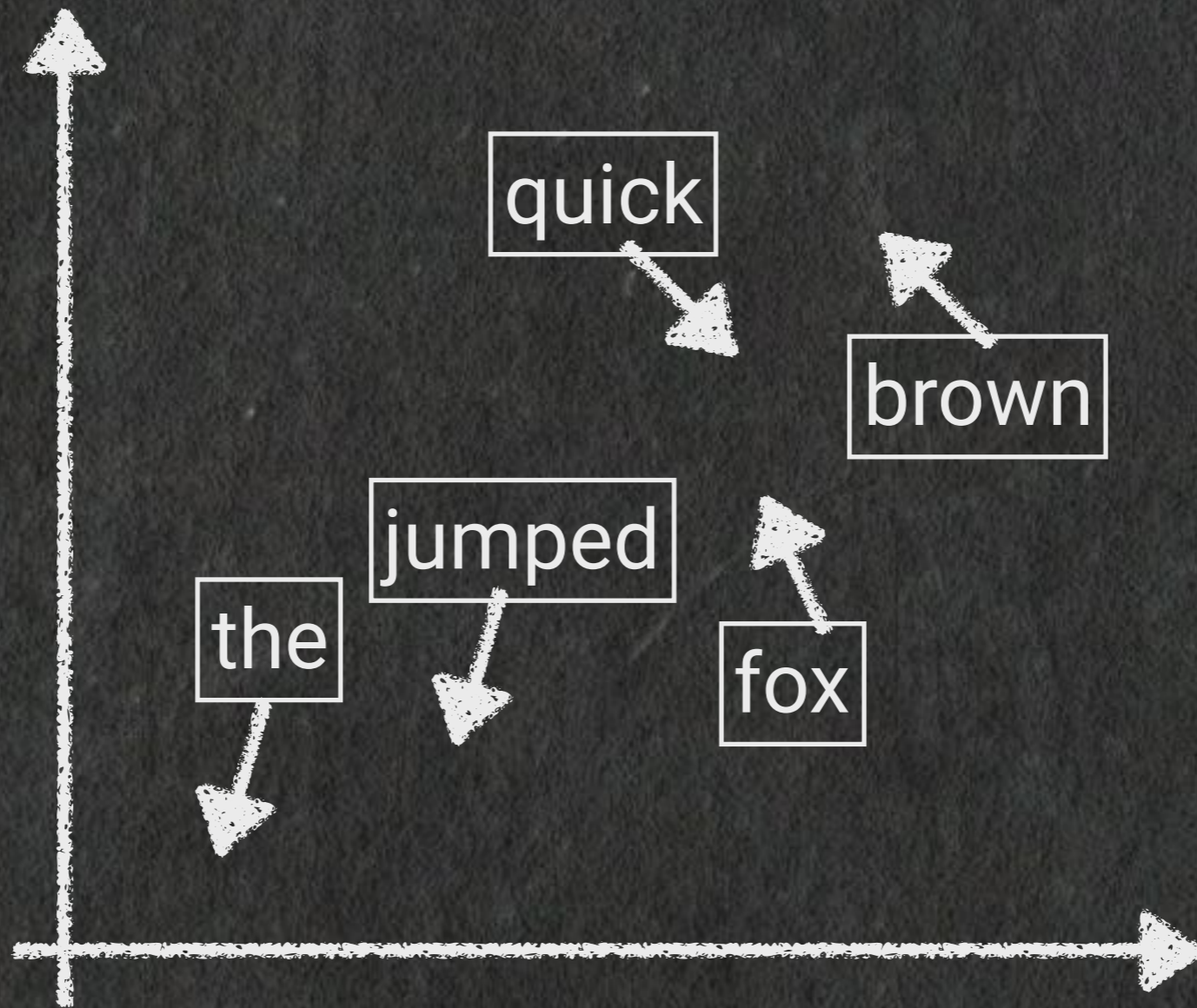
the **quick** brown fox jumped over the lazy dog



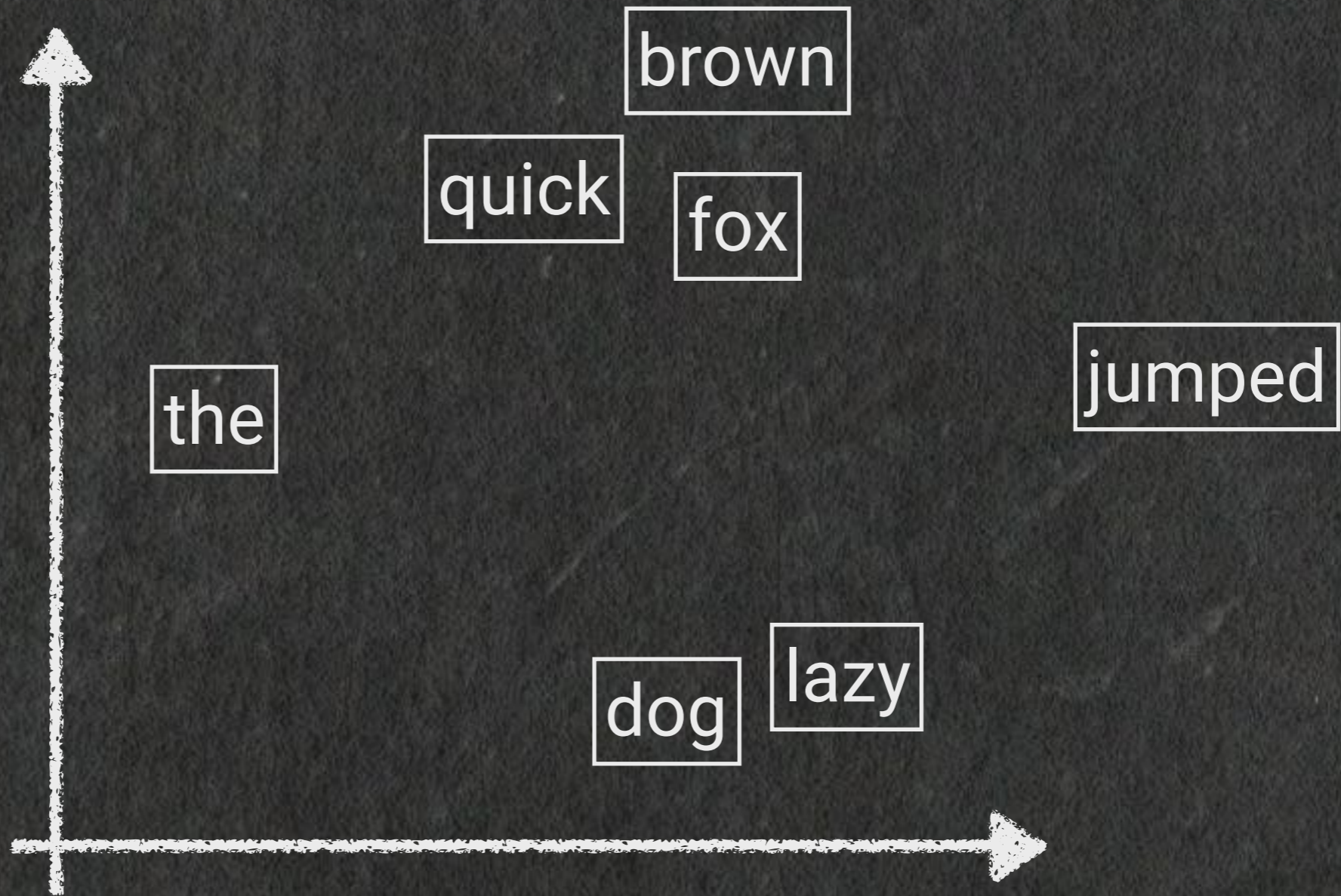
the quick **brown** fox jumped over the lazy dog



~~the quick~~ **brown** ~~fox~~ jumped over the lazy dog
? ?



the quick brown fox jumped over the lazy dog









One size fits all

Criteo

120

billions events

Criteo

Wikipedia

120

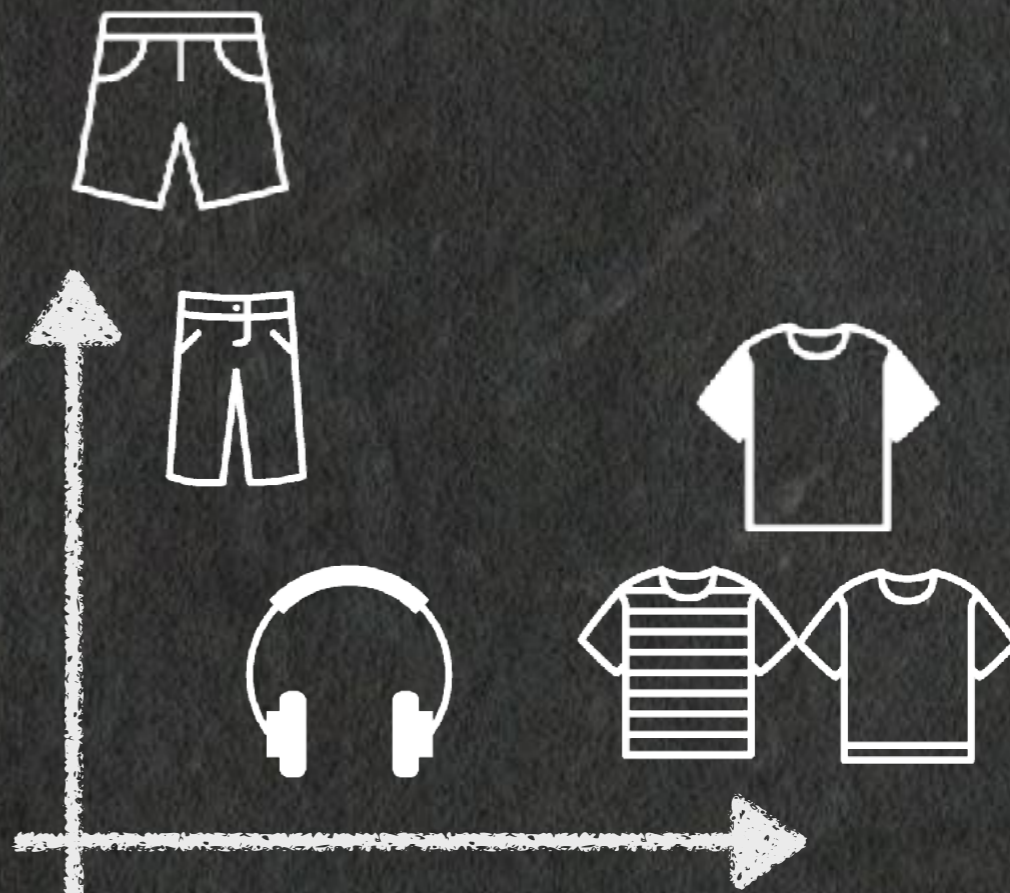
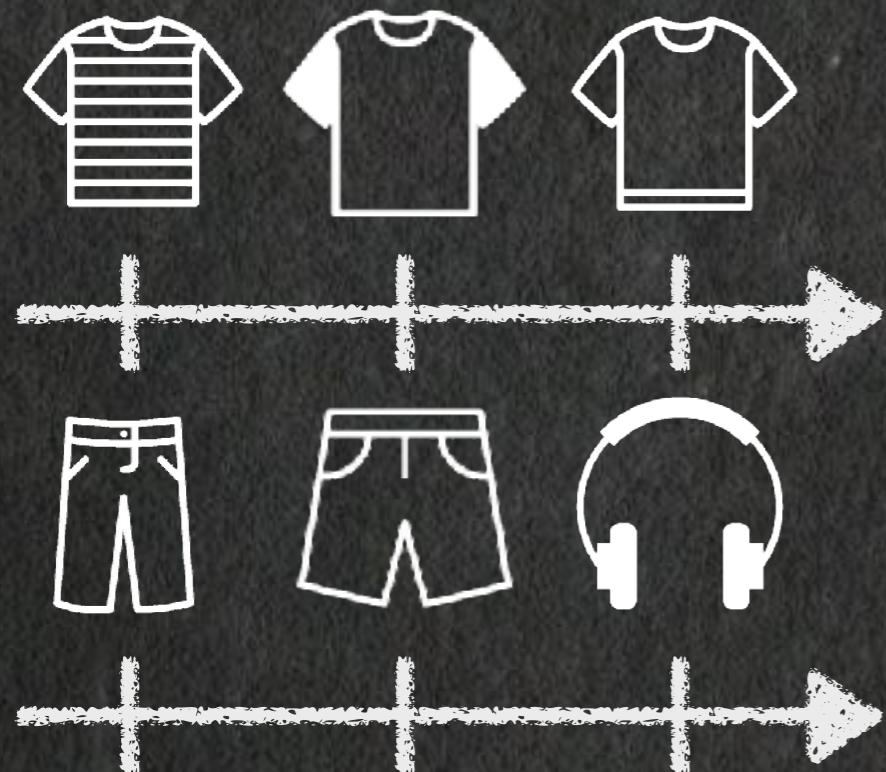
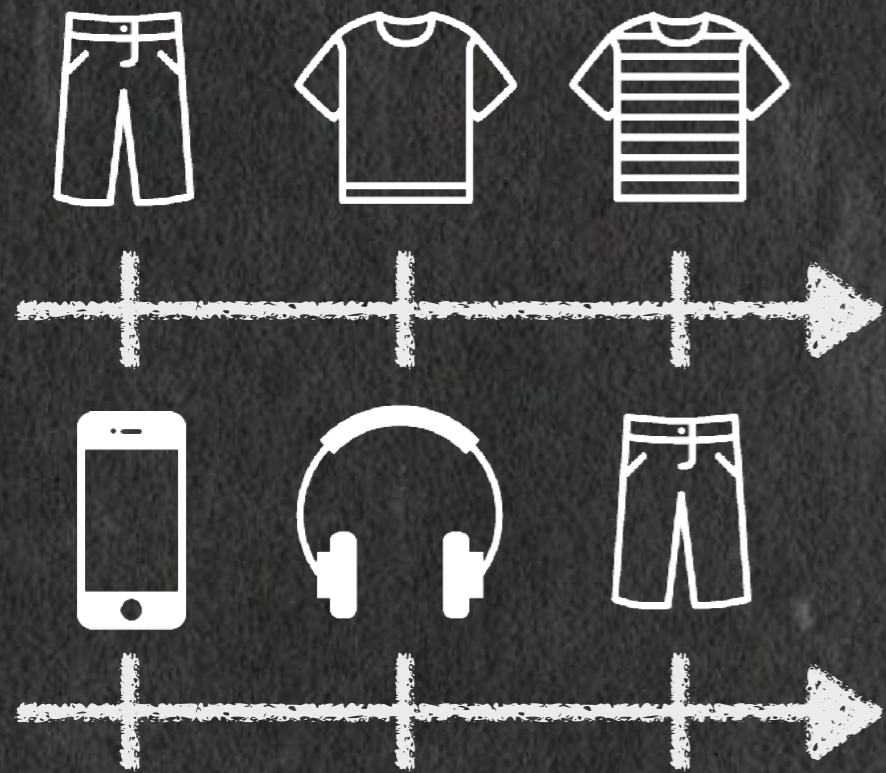
billions events

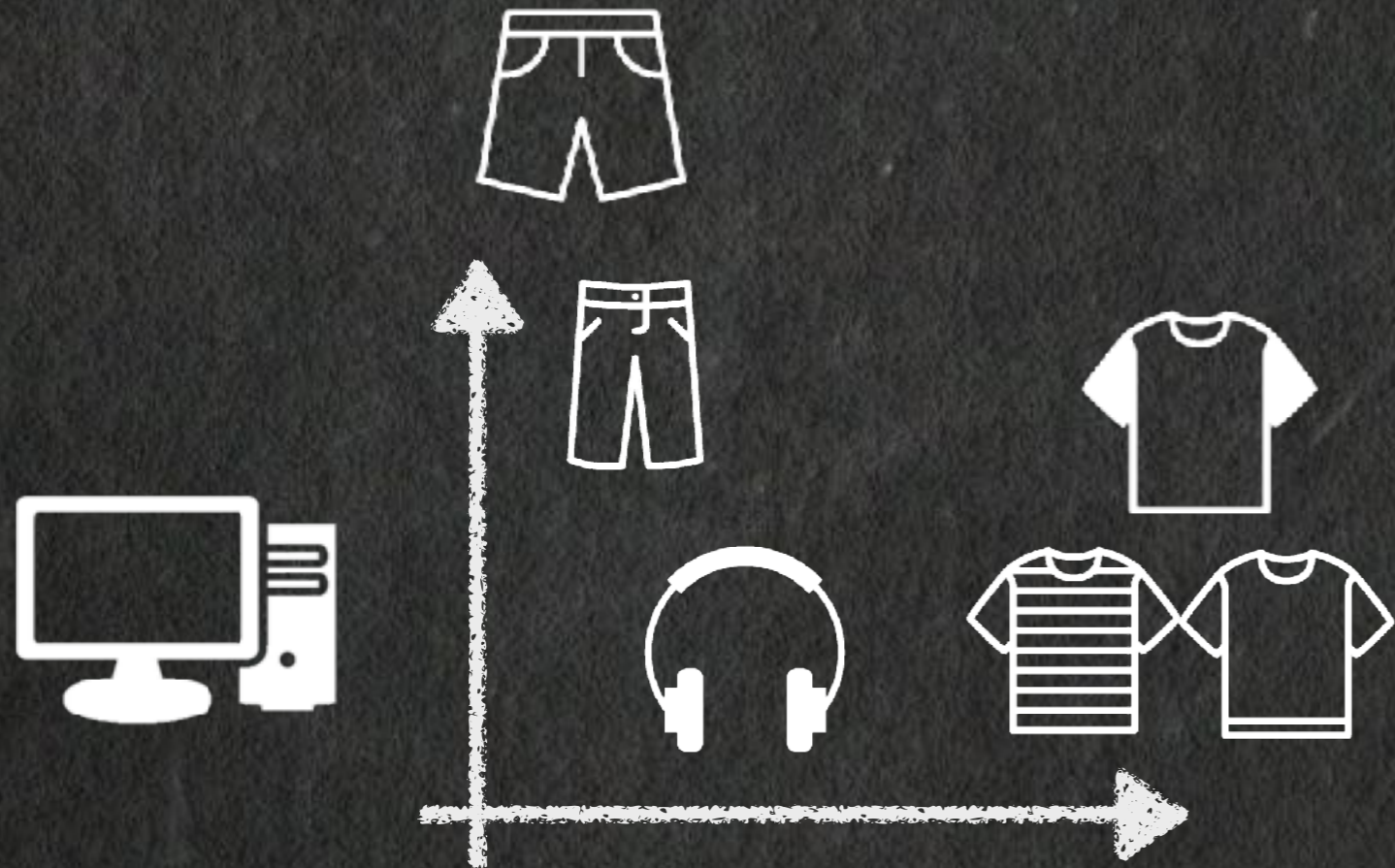
30

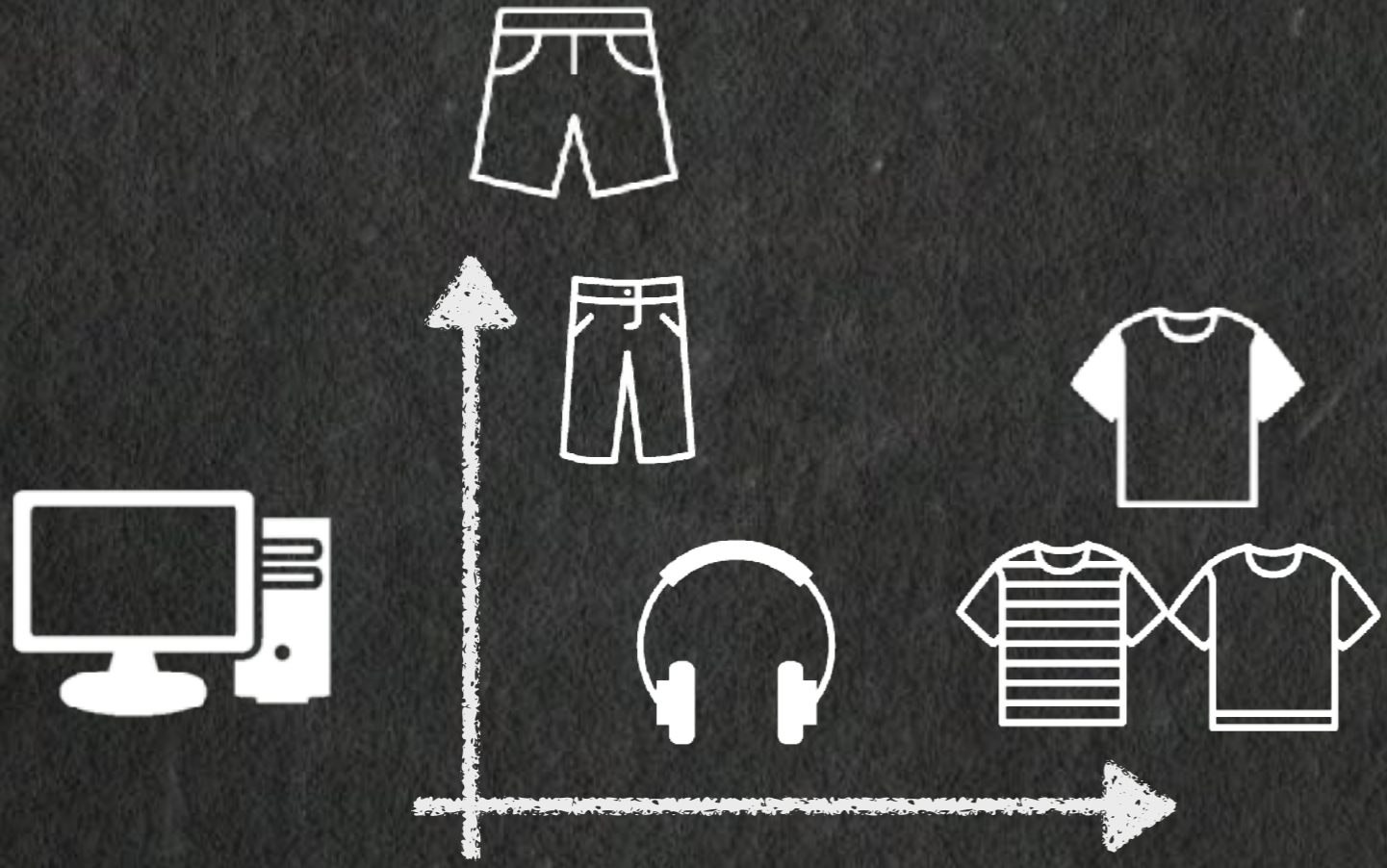
billions words





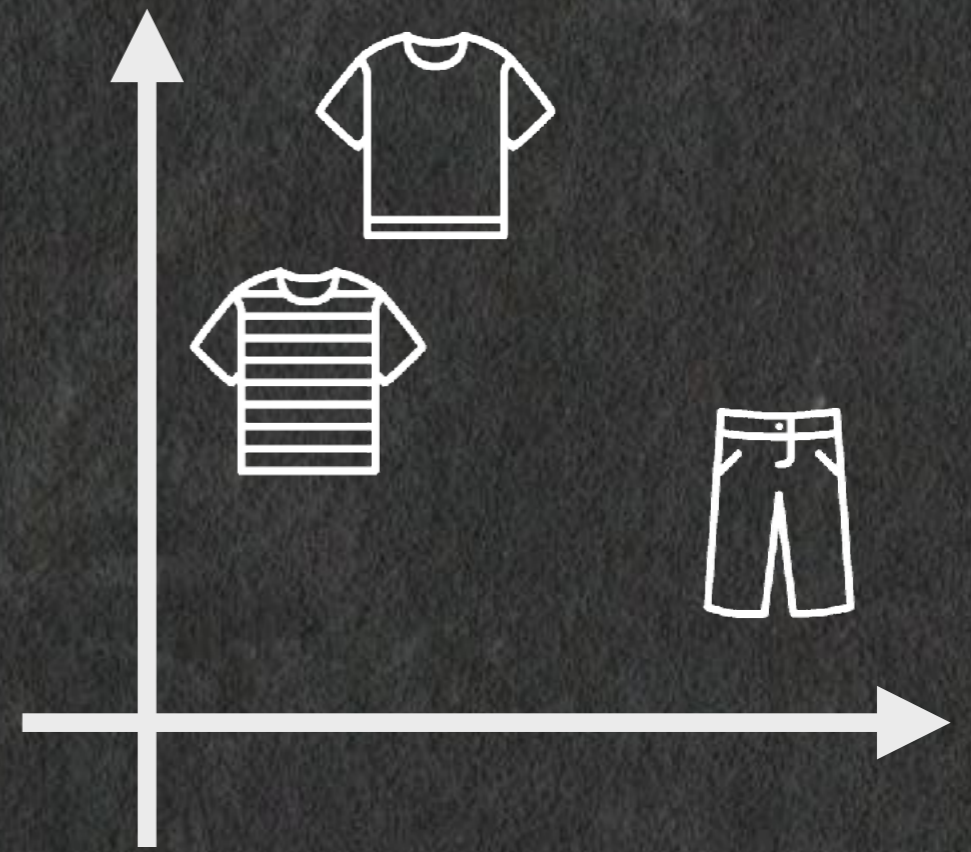


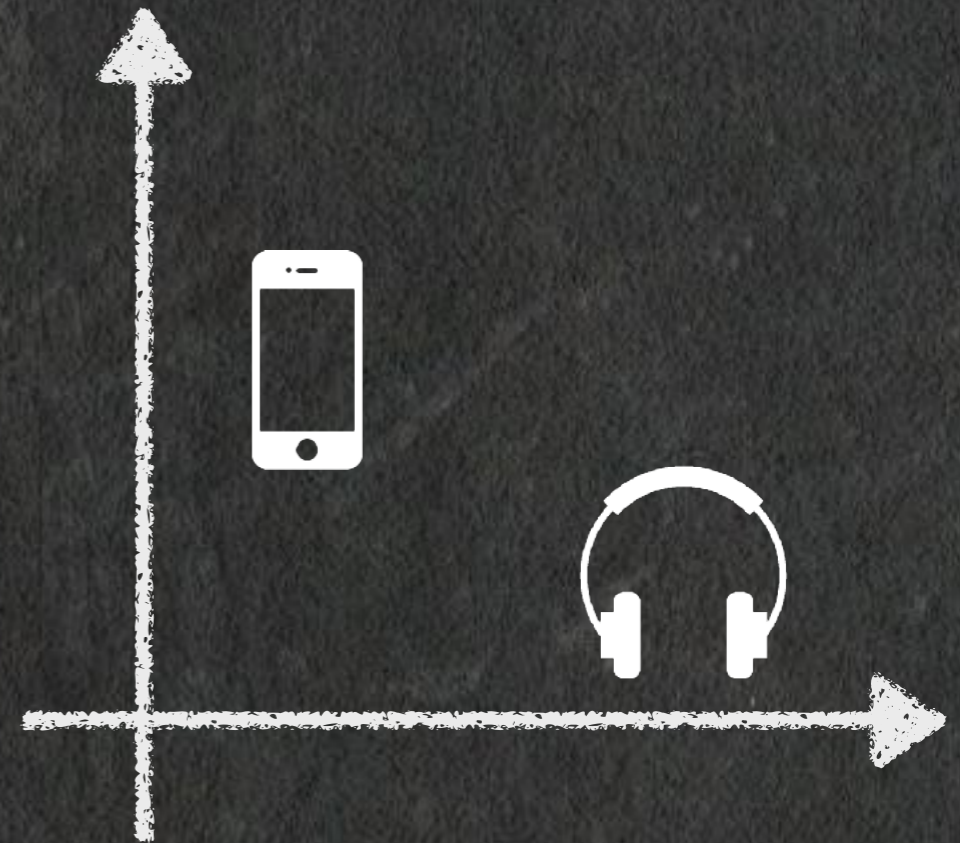
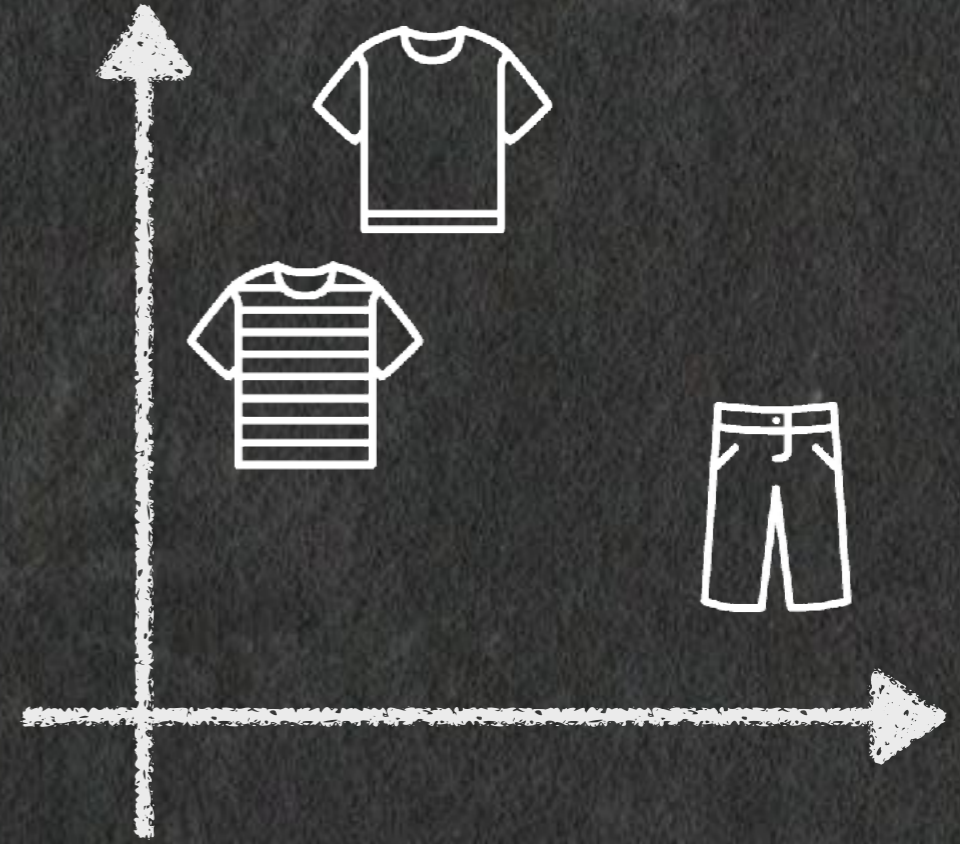












Metrics?
What for?



Lh

measure
probabilities

close to what
we optimize

L1h

measure
probabilities

close to what
we optimize

Recall@K

measure
ranking

close to final
goal

- Number of partitions
- Number of epochs
- Learning rate
- Embedding dimension



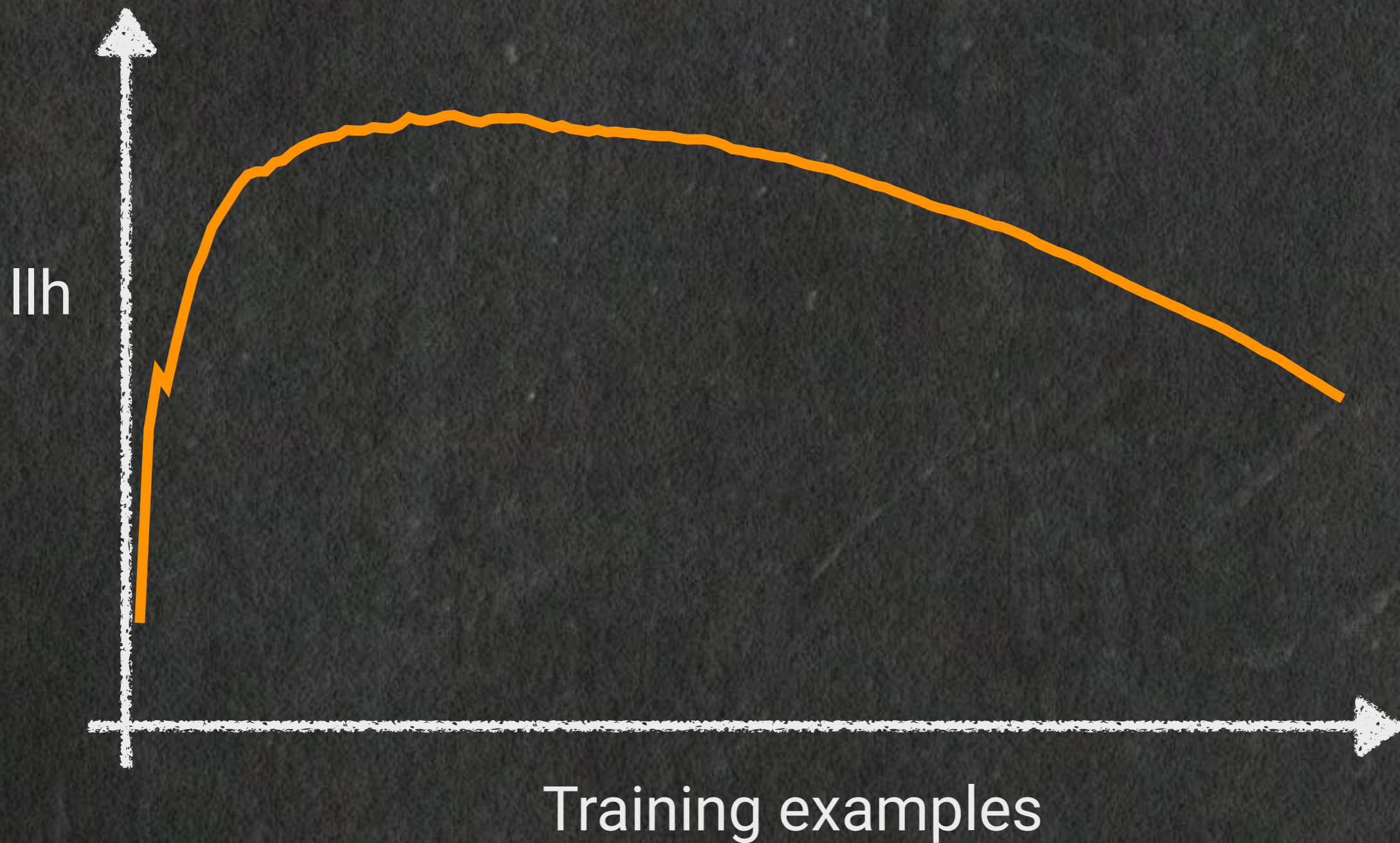
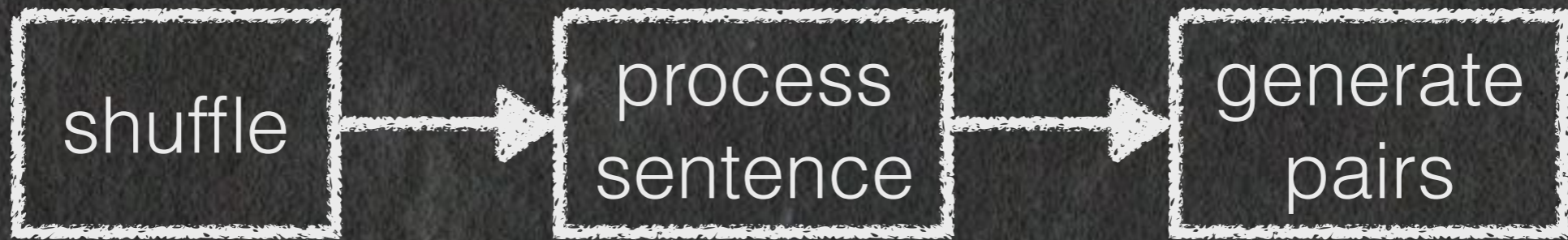
What about **clicks**?

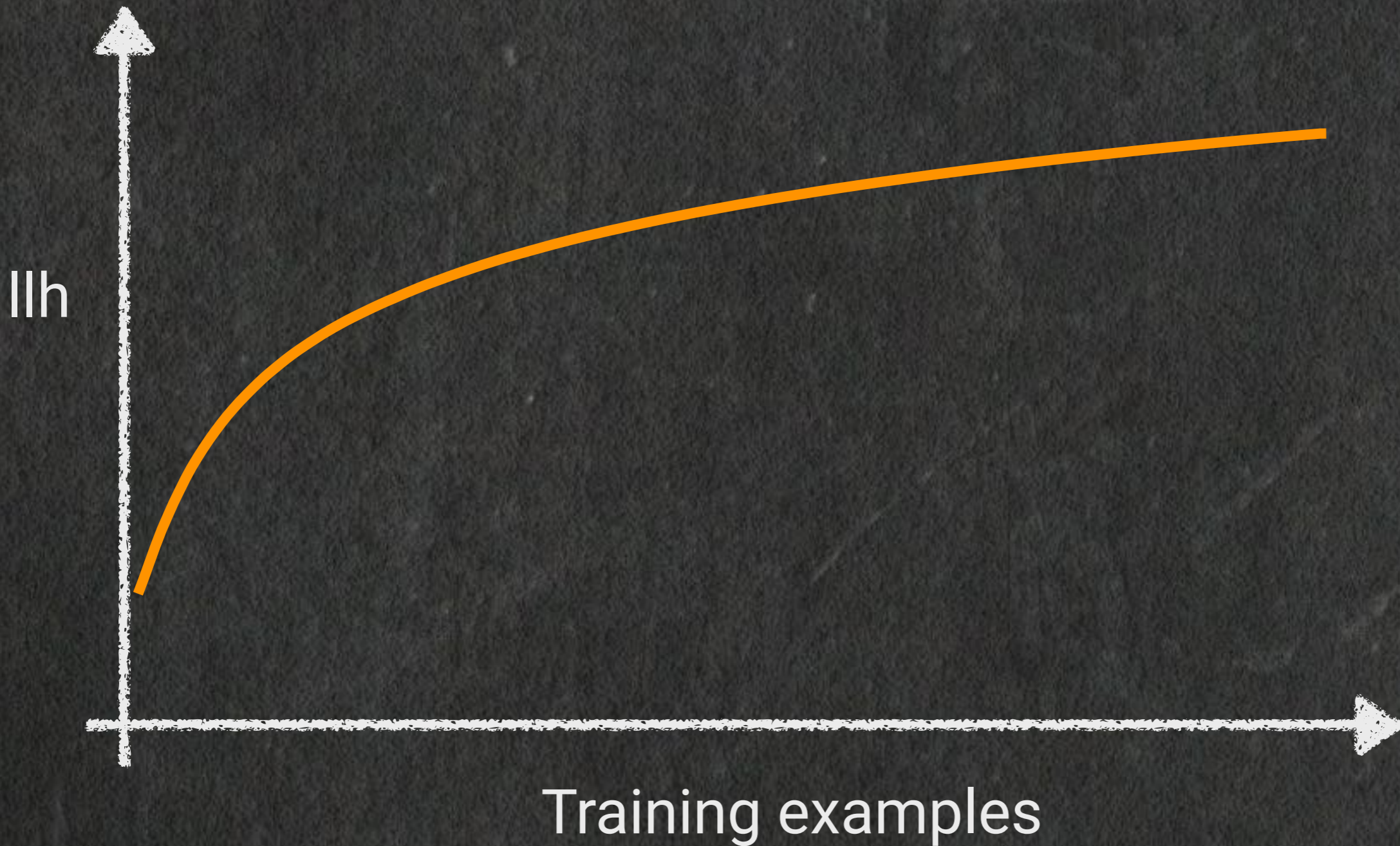
Data

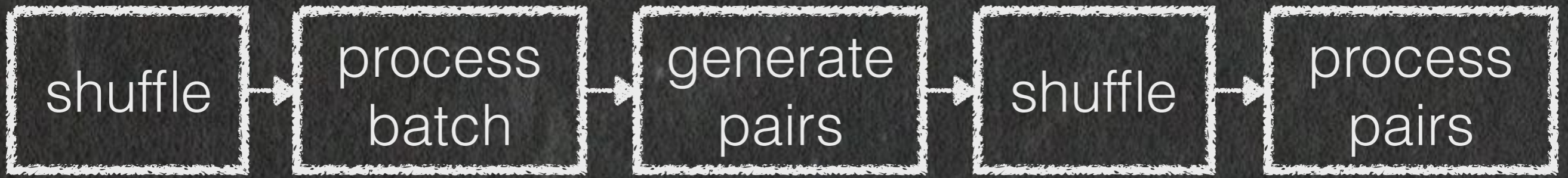
betrayed us!

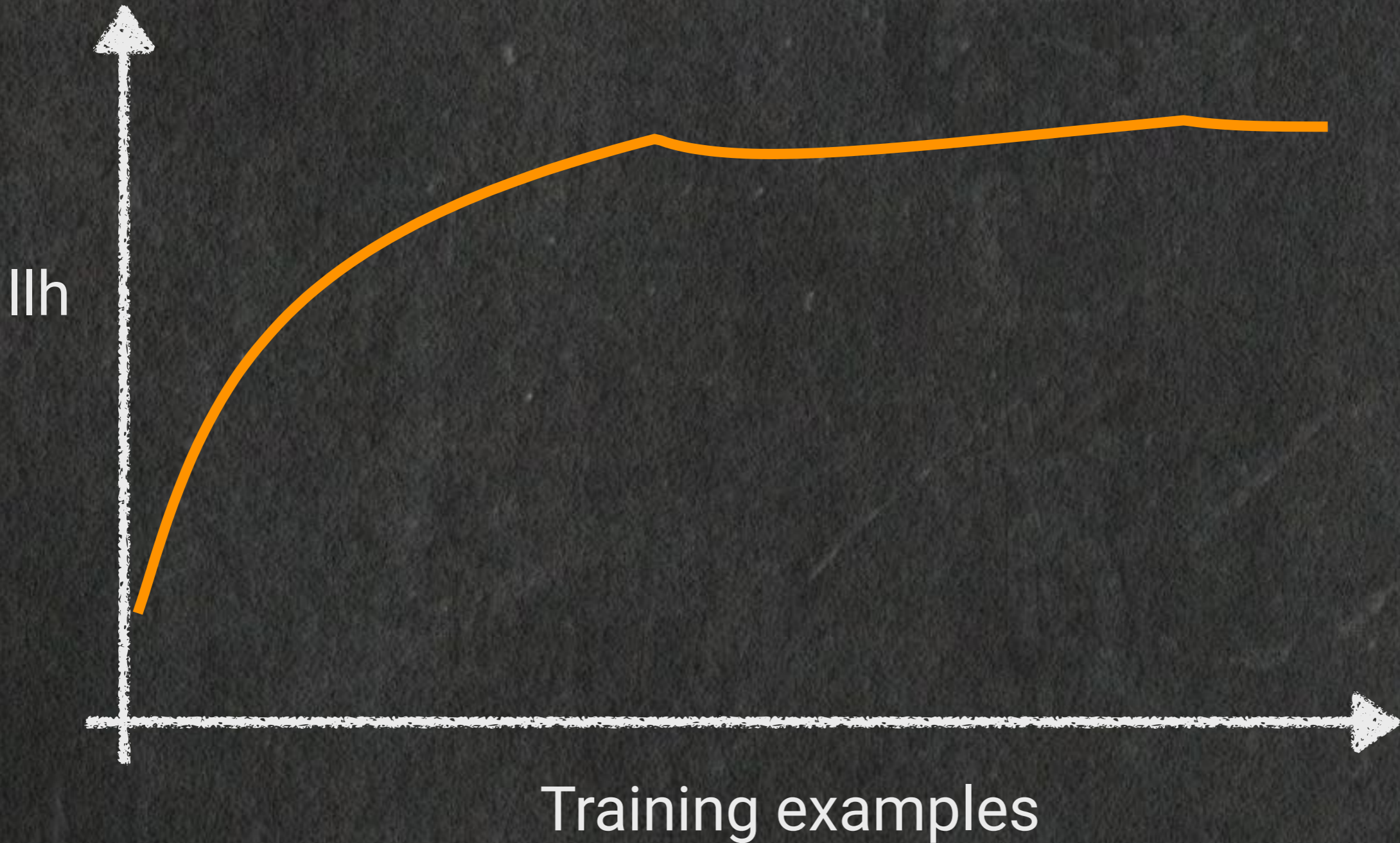
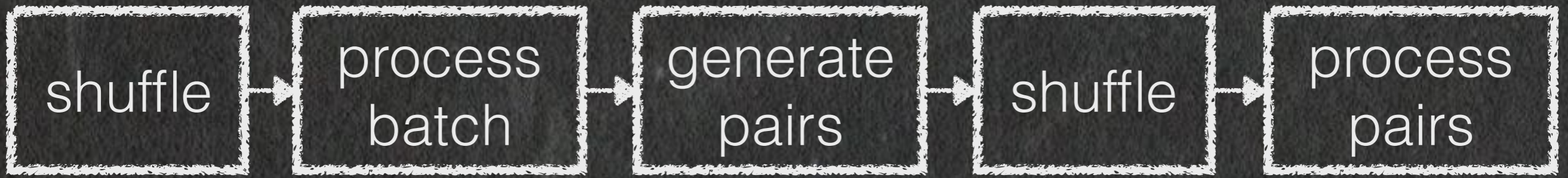












Embrace **trade-offs**

Embrace **trade-offs**

Measure everything

Embrace **trade-offs**

Measure everything

Play with your data

We're hiring!

@simondolle

<http://labs.criteo.com/>

Credits: Ralf Schmitzer, Edward Boatman, Bernar Novalyi,
Royyan Wijaya, art shop, iconsphere

criteo.