



LMs Inside Out

Sasha Rush

youtube.com/@srush_nlp





LMs Inside Out

Sasha Rush

youtube.com/@srush_nlp





LMs Inside Out

Sasha Rush

youtube.com/@srush_nlp

The logo for the Summit on Responsible Decentralized Intelligence. It features a silhouette of a clock tower against a sunset or sunrise sky. Text on the logo includes: "Summit on Responsible Decentralized Intelligence", "--- Future of Decentralization and AI", "August 6, 2024, New York City", "rdi.berkeley.edu", "Hosted By Berkeley Decentralized Intelligence", and a quote at the bottom: "the work that's been going on in these areas kind of pulling everything together."

Collaborators



Jack Morris



Volodymyr Kuleshov



Justin Chiu



Vitaly Shmatikov



Wenting Zhao

This one,

Collaborators



Jack Morris



Volodymyr Kuleshov



Justin Chiu



Vitaly Shmatikov

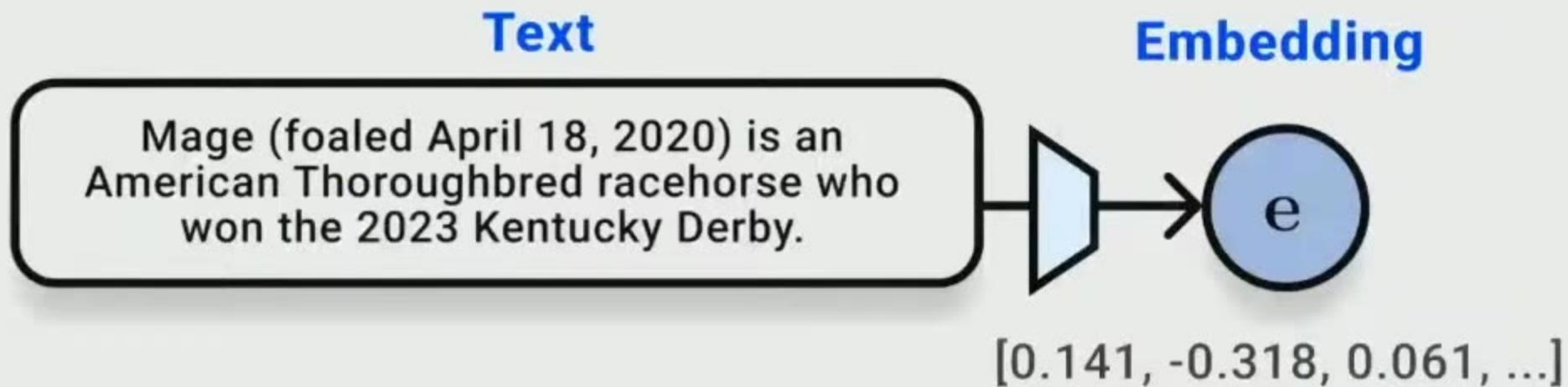


Wenting Zhao

This one. Yeah, beautiful. Okay. Great. Cool



Text embeddings

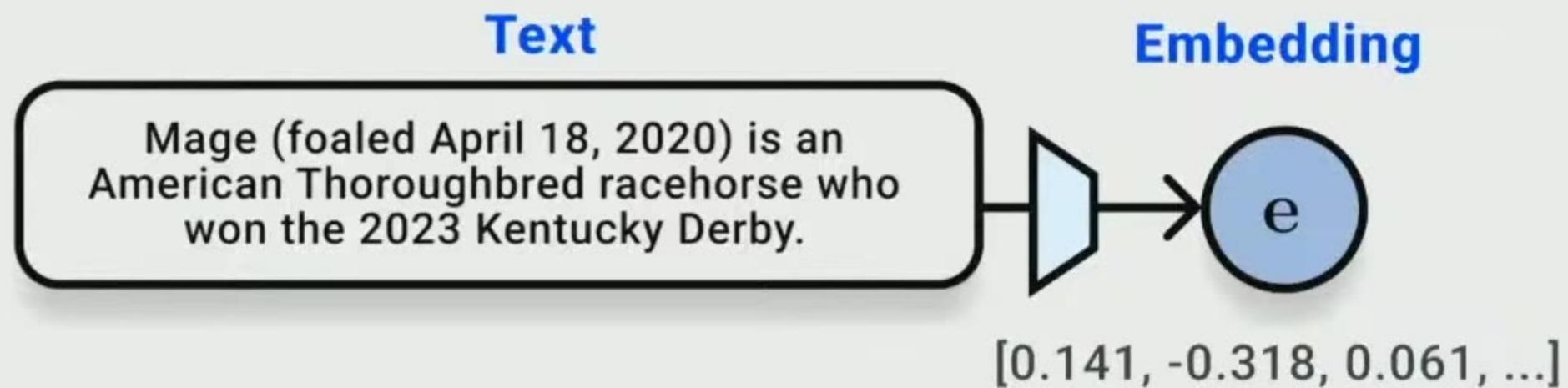


This one. Yeah, beautiful. Okay, great. Cool. So this



Verizon Event ...

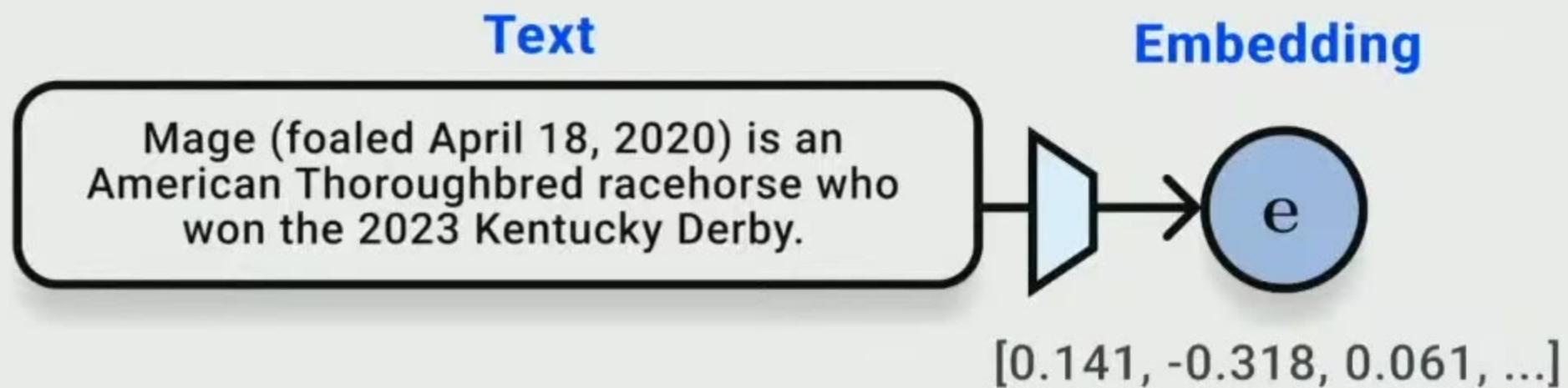
Text embeddings



Okay, so in this paper, we're going to



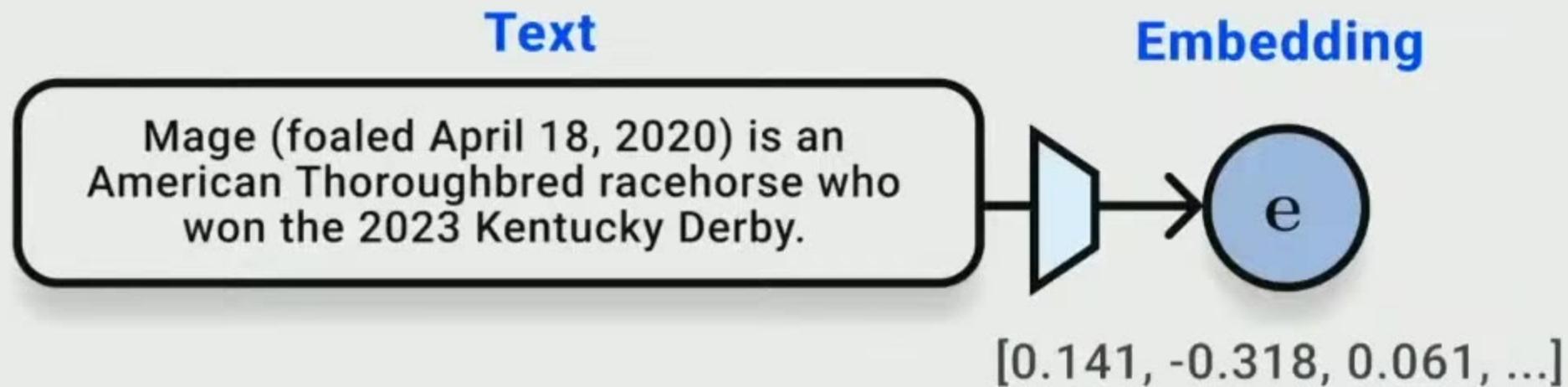
Text embeddings



so in this paper, we're going to be talking about text embeddings.



Text embeddings



Text. Embeddings, I think, are probably the most widely

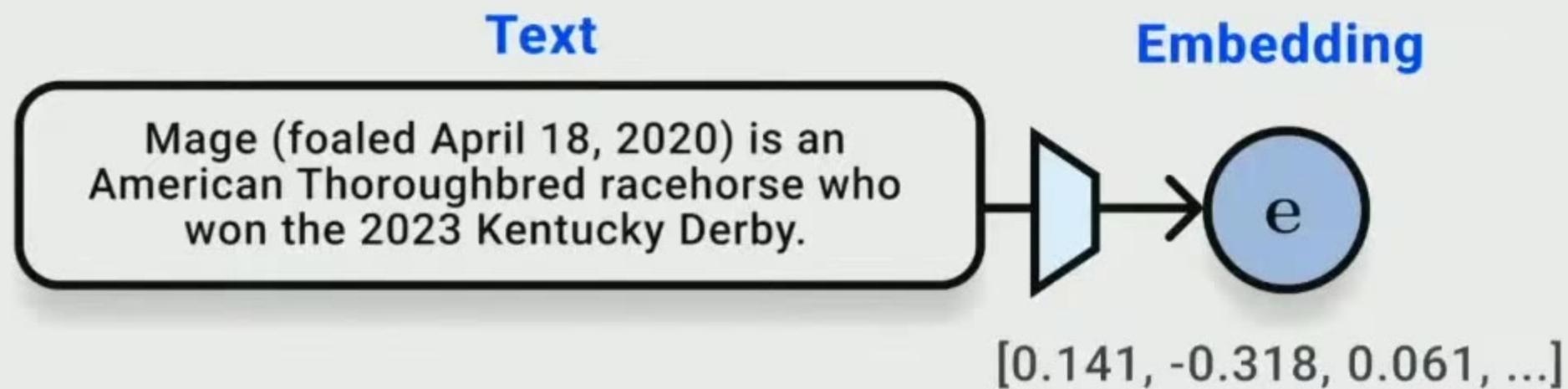


Text embeddings





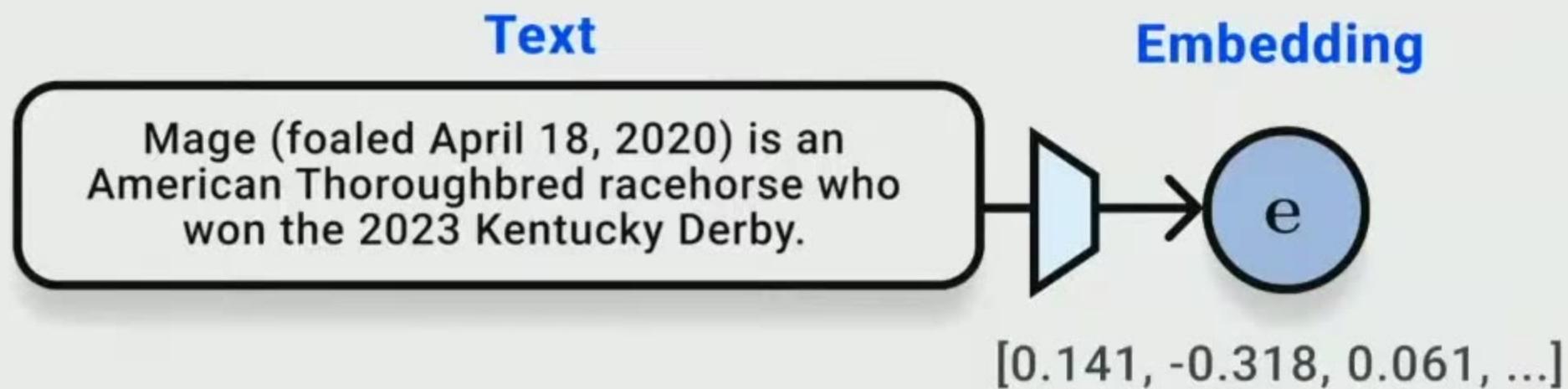
Text embeddings



People use them for everything, from classification to



Text embeddings



retrieval. All sorts of applications of that form.



Vector databases are *huge* right now

Startups

Pinecone drops \$100M in funding on \$750M valuation, as vector database demand grows

AI

Qdrant, an open source vector database, develops AI data

Paul Sawers @psawers

Premium HOME > TECH

Vector database Chroma scored \$18 million valuation. Here's why its technology is fueling generative AI startups.

Stephanie Palazzolo Apr 6, 2023, 8:00 AM EDT

Weaviate Raises \$50 Million Series B Funding
Meet Soaring Demand for AI Native Vector Database Technology

NEWS PROVIDED BY
[Weaviate](#) →
21 Apr, 2023, 08:00 ET

FORBES > INNOVATION > CLOUD

The Rise Of Vector Databases

Company's open source vector database and new cloud service

So couple of years ago there were about

Adrian Bridgwater Senior Contributor ⓘ

Berkeley
UNIVERSITY OF CALIFORNIA
Powered by Zoom



Vector databases are *huge* right now

Startups

Pinecone drops \$100M in funding on \$750M valuation, as vector database demand grows

AI

Qdrant, an open source vector database, develops AI tools for data

Paul Sawers @psawers

Premium HOME > TECH

Vector database Chroma scored \$18 million valuation. Here's why its technology is fueling generative AI startups.

Stephanie Palazzolo Apr 6, 2023, 8:00 AM EDT

Weaviate Raises \$50 Million Series B Funding
Meet Soaring Demand for AI Native Vector Database Technology

NEWS PROVIDED BY

Weaviate →

21 Apr, 2023, 08:00 ET

FORBES > INNOVATION > CLOUD

The Rise Of Vector Databases

Company's open source vector database and new cloud service are changing how companies that started that kind of offered a service

Adrian Bridgwater Senior Contributor ⓘ

Berkeley
UNIVERSITY OF CALIFORNIA
Powered by Zoom



Vector databases are *huge* right now

Startups

Pinecone drops \$100M in funding on \$750M valuation, as vector database demand grows

AI

Qdrant, an open source vector database, develops AI data

Paul Sawers @psawers

Premium HOME > TECH

Vector database Chroma scored \$18 million valuation. Here's why its technology is fueling generative AI startups.

Stephanie Palazzolo Apr 6, 2023, 8:00 AM EDT

Weaviate Raises \$50 Million Series B Funding
Meet Soaring Demand for AI Native Vector Database Technology

NEWS PROVIDED BY
[Weaviate](#) →
21 Apr, 2023, 08:00 ET

Company's open source vector database and new cloud service aim to make AI more accessible

FORBES > INNOVATION > CLOUD

The Rise Of Vector Databases

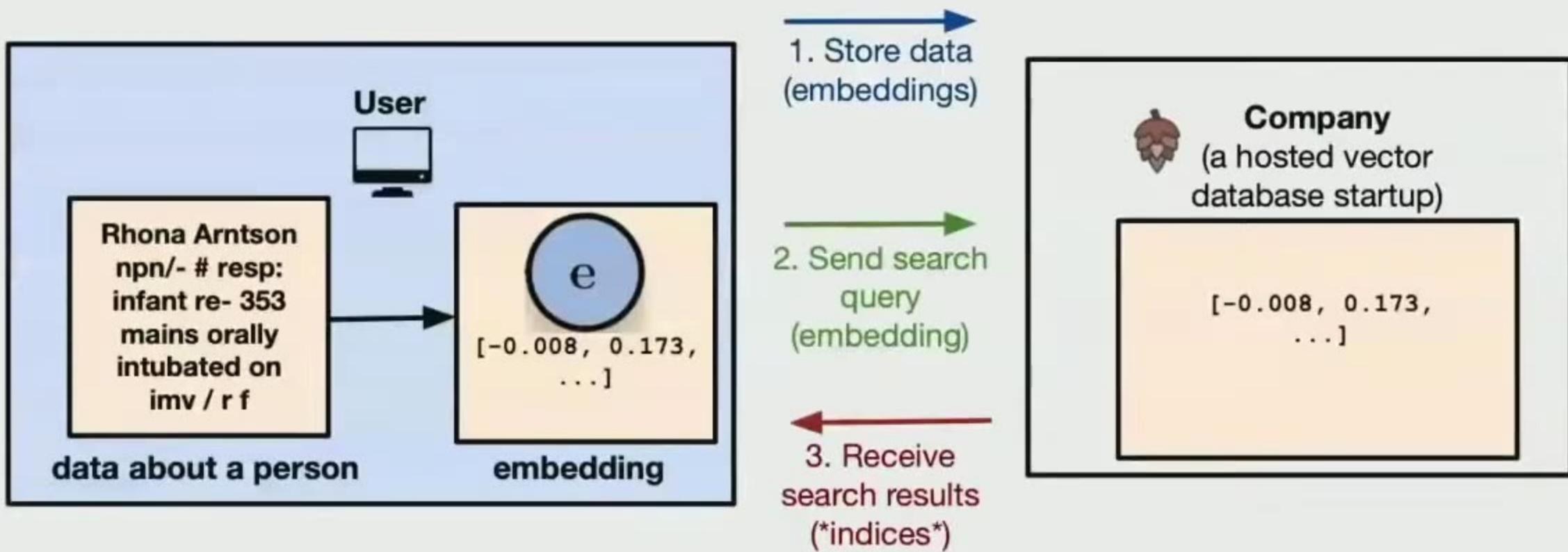
They were tossed into a large

Adrian Bridgwater Senior Contributor ⓘ

Berkeley
UNIVERSITY OF CALIFORNIA
Powered by Zoom



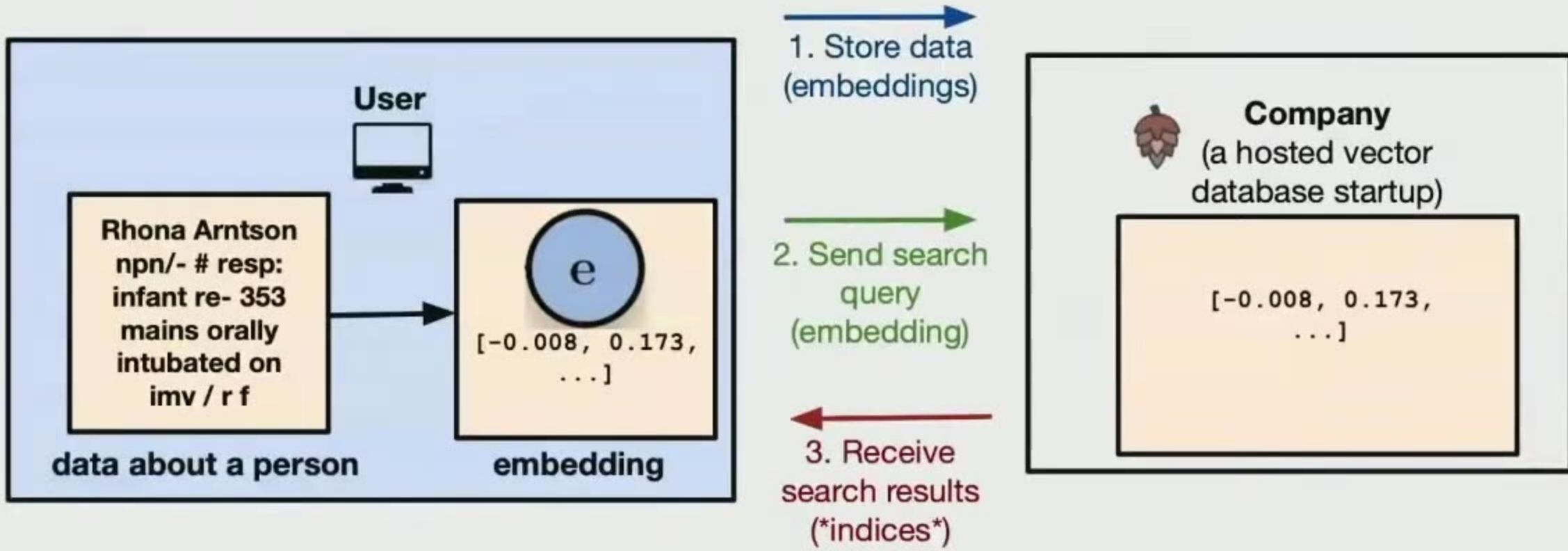
Threat model



They were tossed into a large database, you'd be able to



Threat model

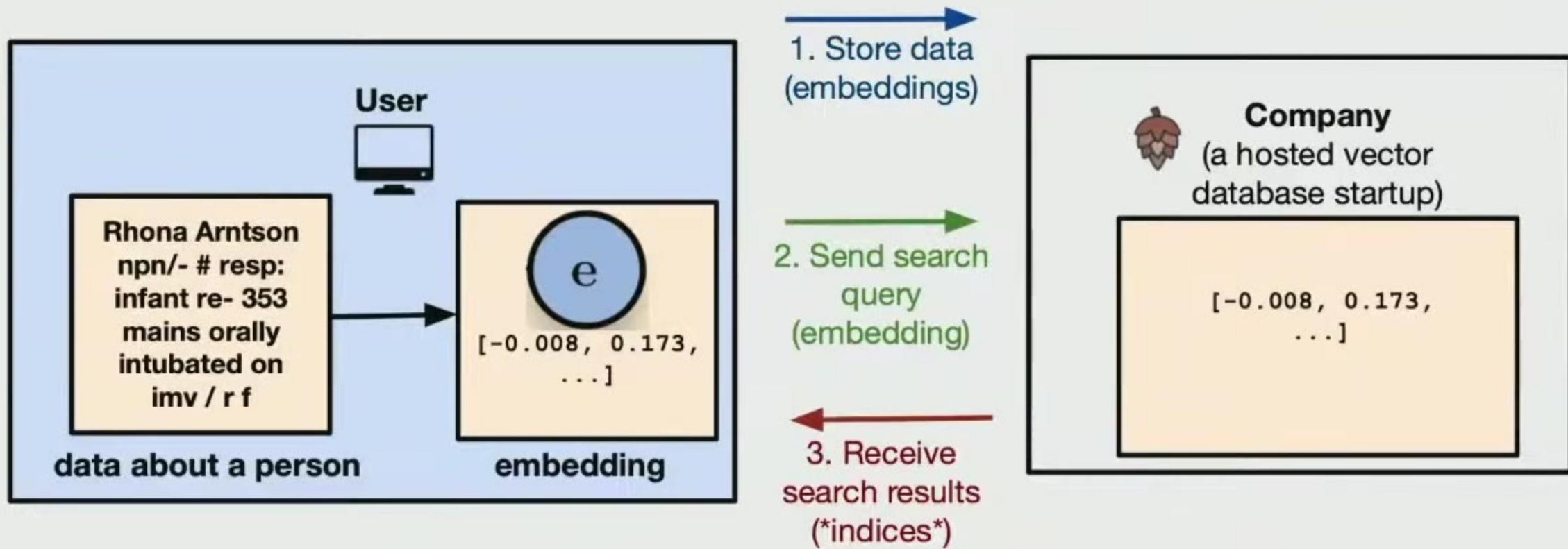


They were tossed into a large database, you'd be able to



Verizon Event ...

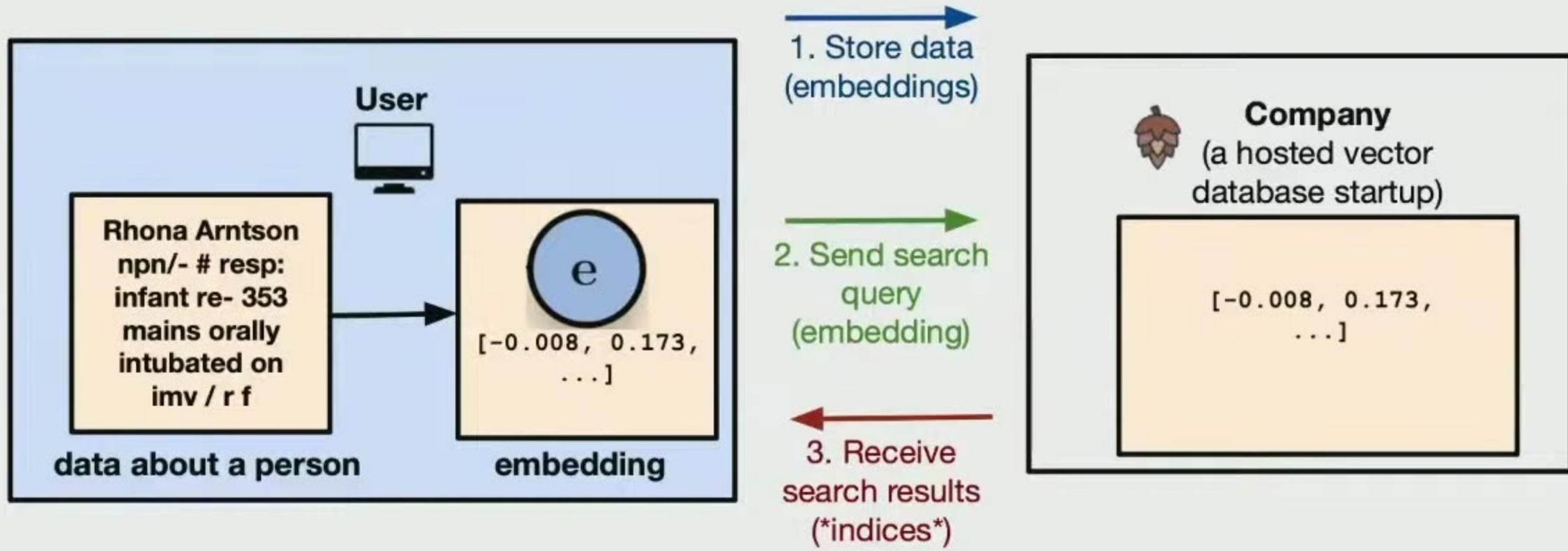
Threat model



We got pretty interested in,



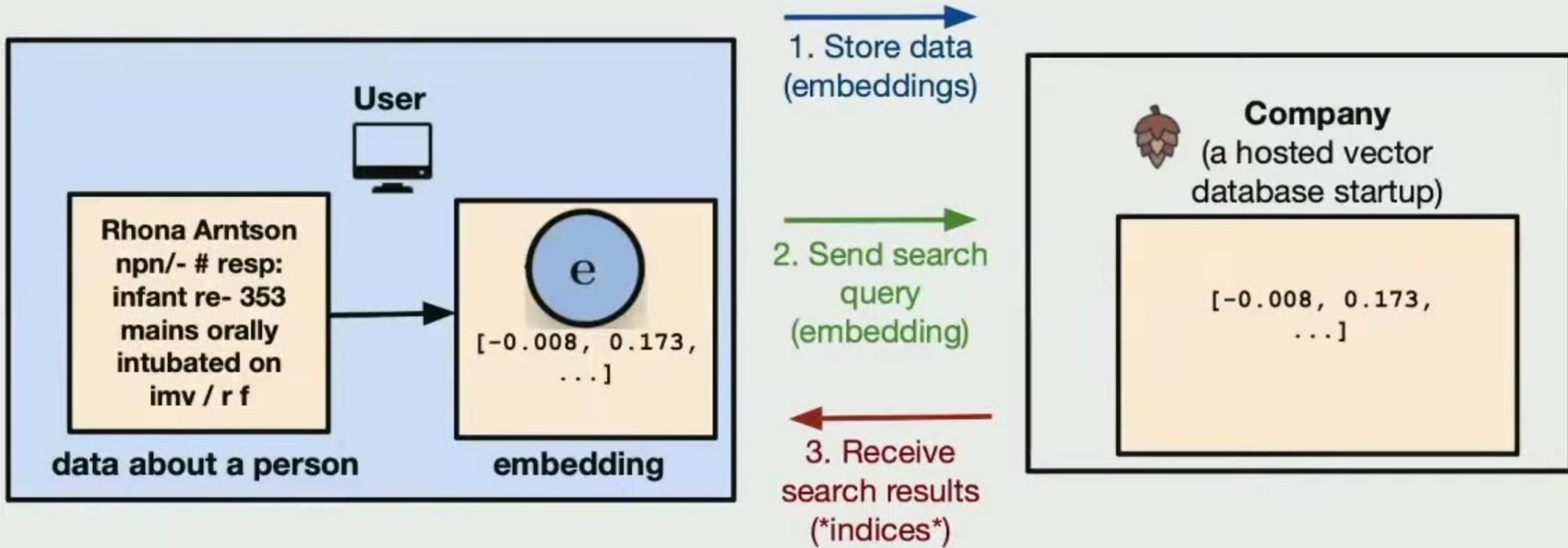
Threat model



for these sort of databases, which is we didn't really



Threat model

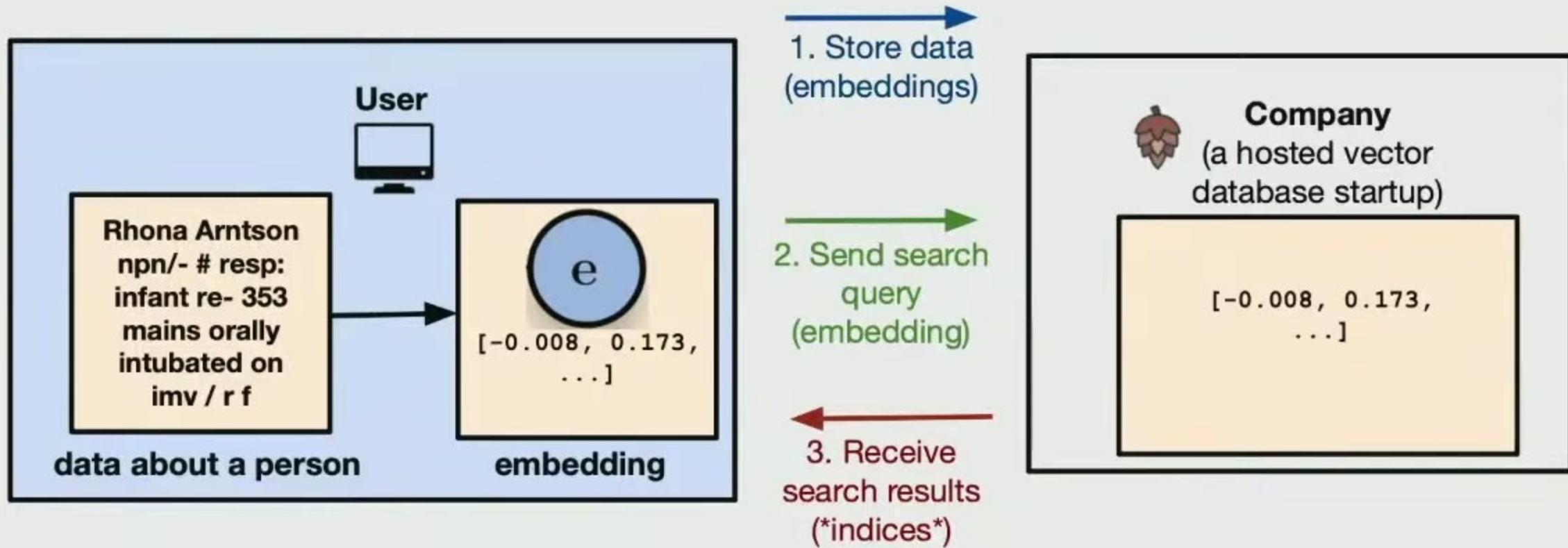


Which is,



Verizon Event ...

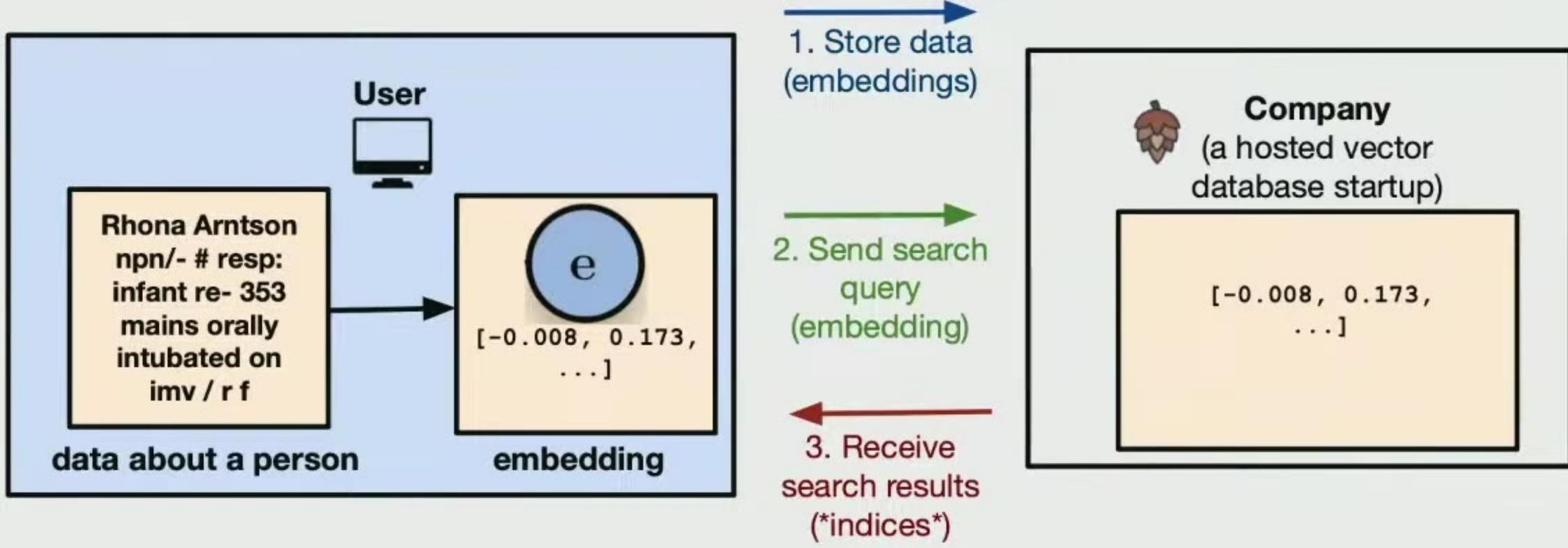
Threat model



is, imagine that a user got access to an embedding database.



Threat model

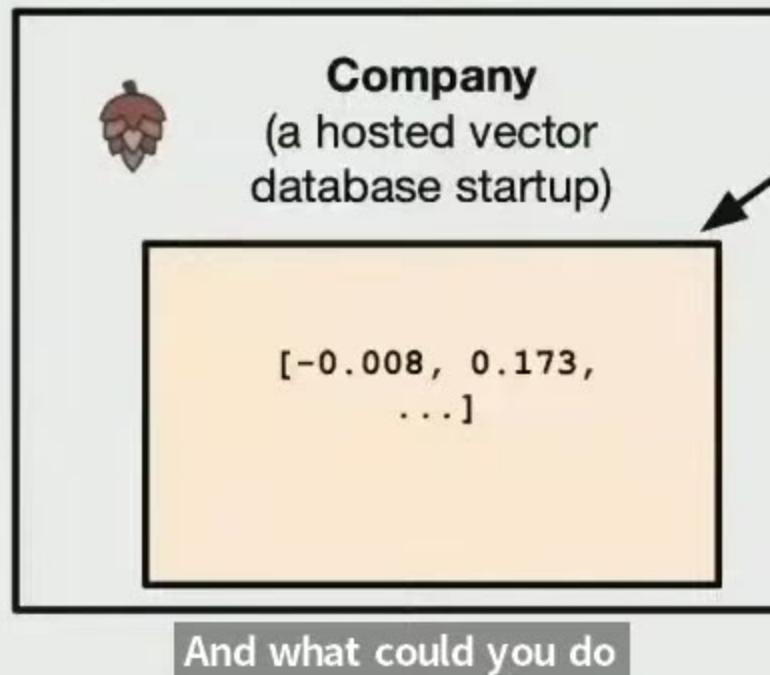


was going to take all their data, send it over to the



Threat model

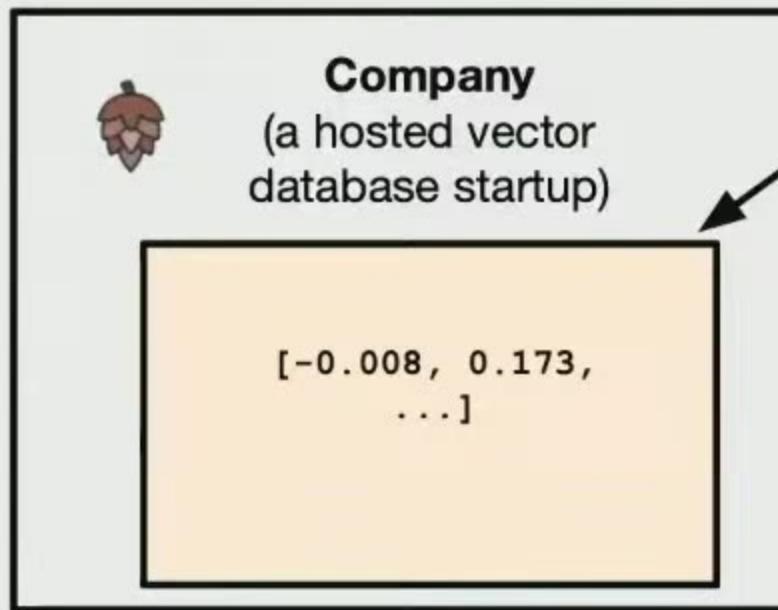
**What can a *bad actor* learn from just looking
embeddings of text?**





Threat model

What can a *bad actor* learn from just looking embeddings of text?



And what could you do in



Inverting text embeddings





Inverting text embeddings





Inverting text embeddings





Why is it hard to invert text embeddings?

- Embedding models are trained to have maximum similarity between two similar pieces of text
- But we still find small numerical differences between different paragraphs' embeddings, even with just one word substituted with a synonym



Why is it hard to invert text embeddings?

- Embedding models are trained to have maximum similarity between two similar pieces of text
- But we still find small numerical differences between different paragraphs' embeddings, even with just one word substituted with a synonym



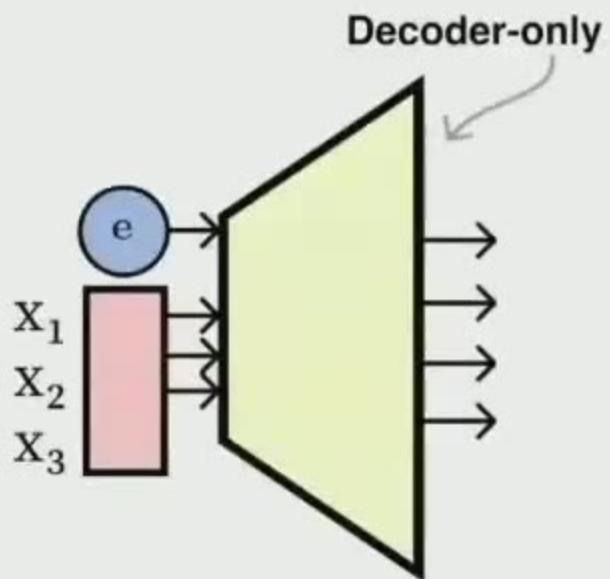
Baseline approach: conditional generation

> How can we generate text from an embedding?



Baseline approach: conditional generation

Prior architecture

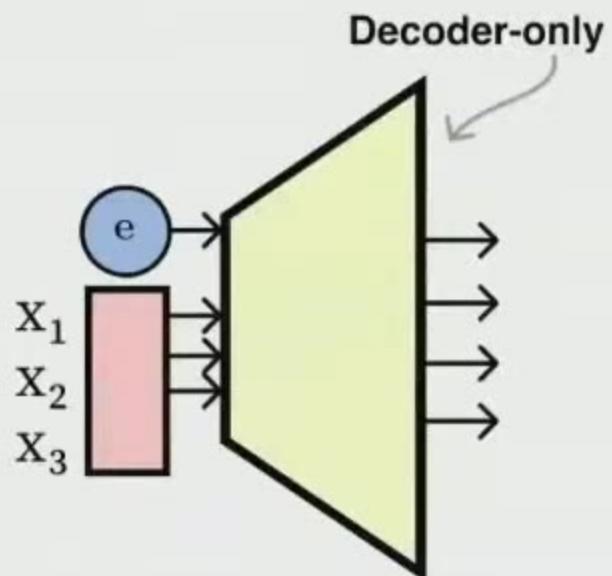


(Adolphs et al., 2022; Li et al., 2023)



Baseline approach: conditional generation

Prior architecture

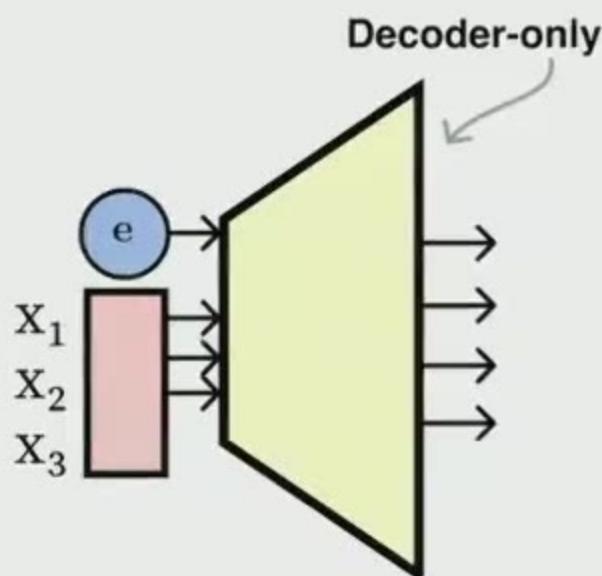


(Adolphs et al., 2022; Li et al., 2023)



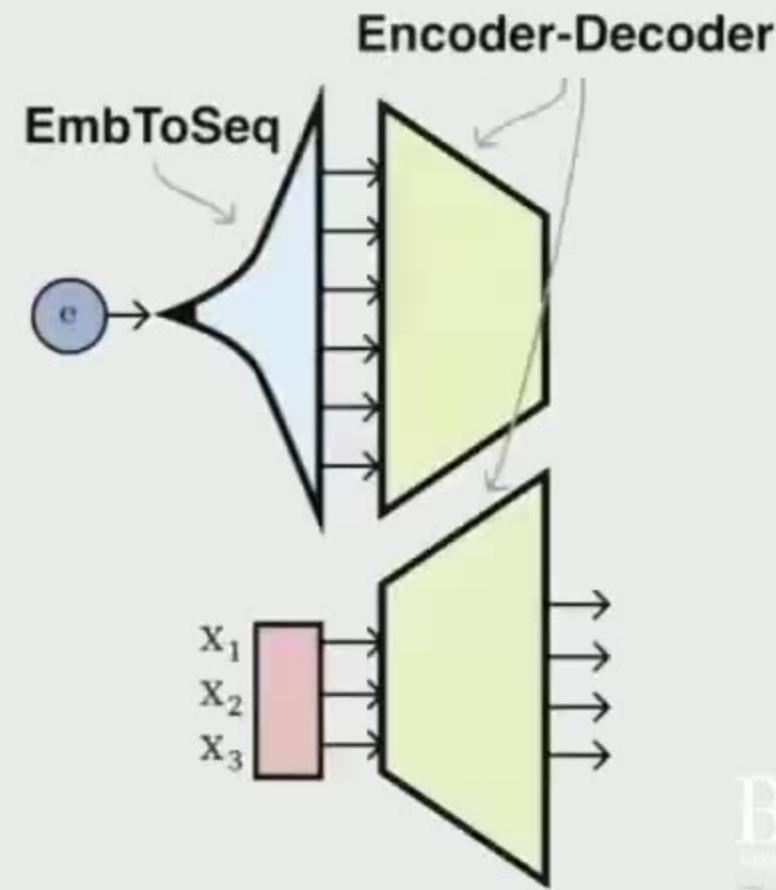
Baseline approach: conditional generation

Prior architecture



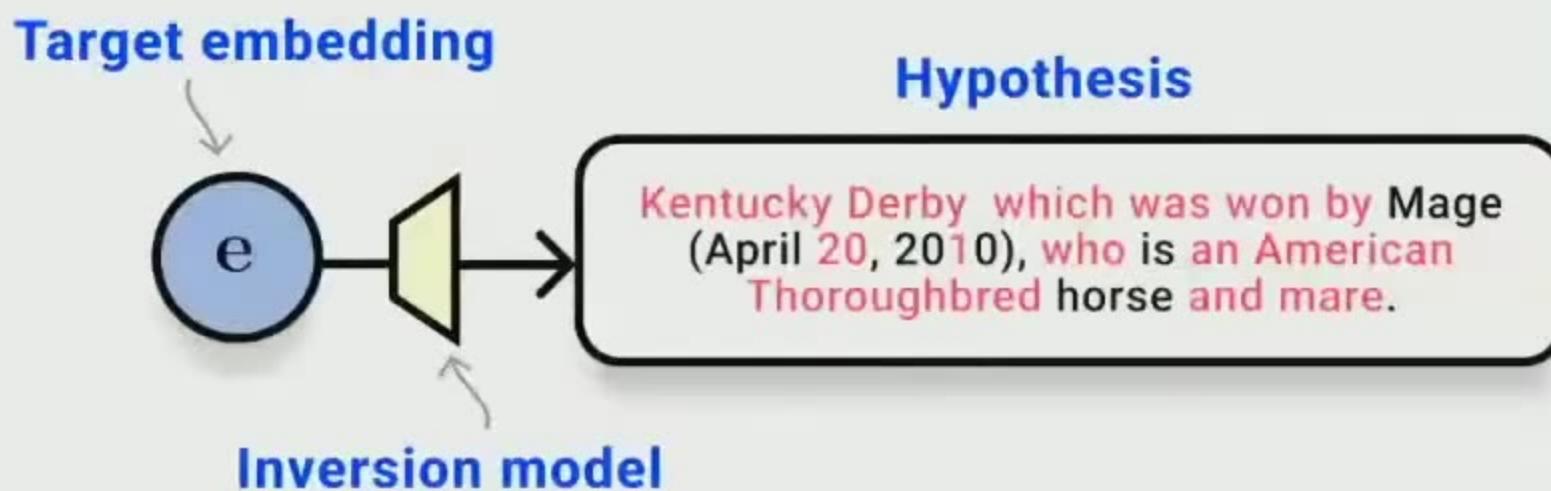
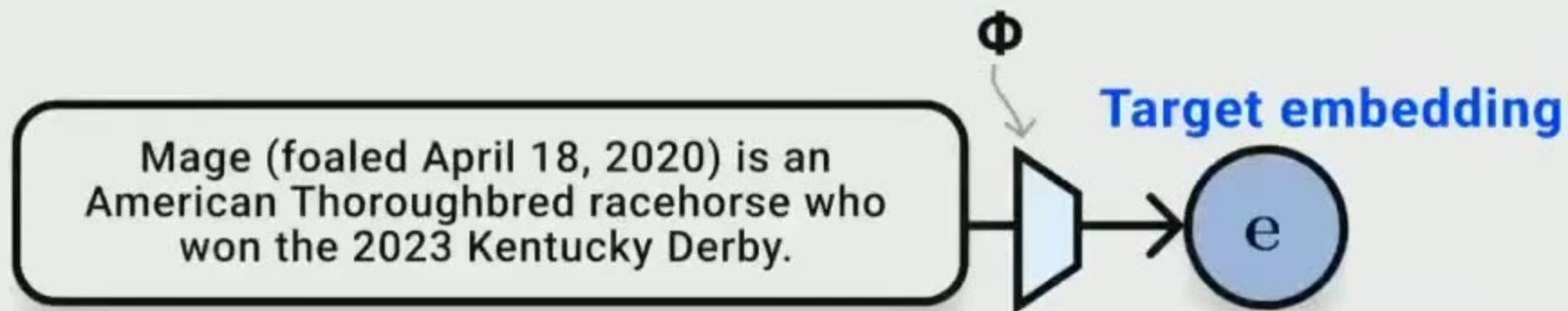
(Adolphs et al., 2022; Li et al., 2023)

Initial architecture (ours)



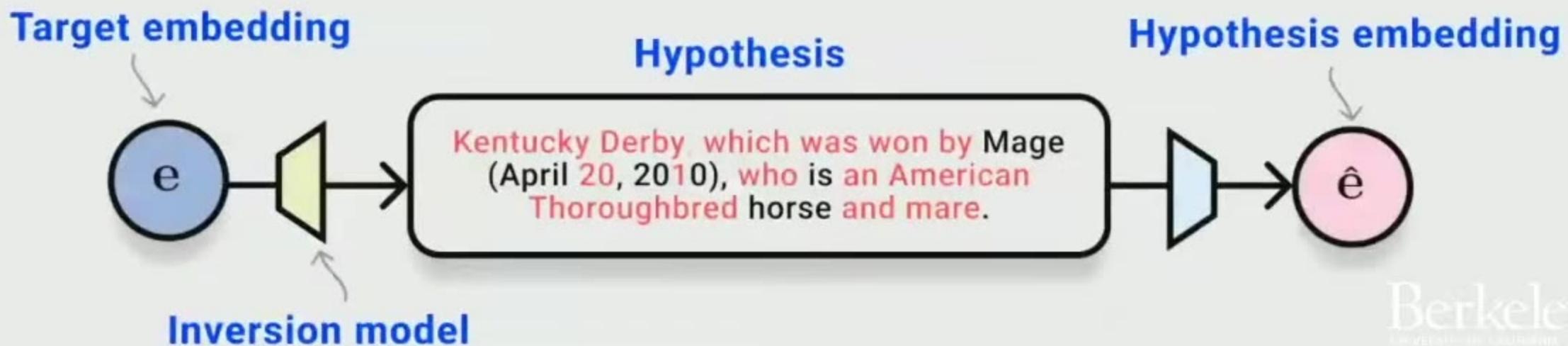
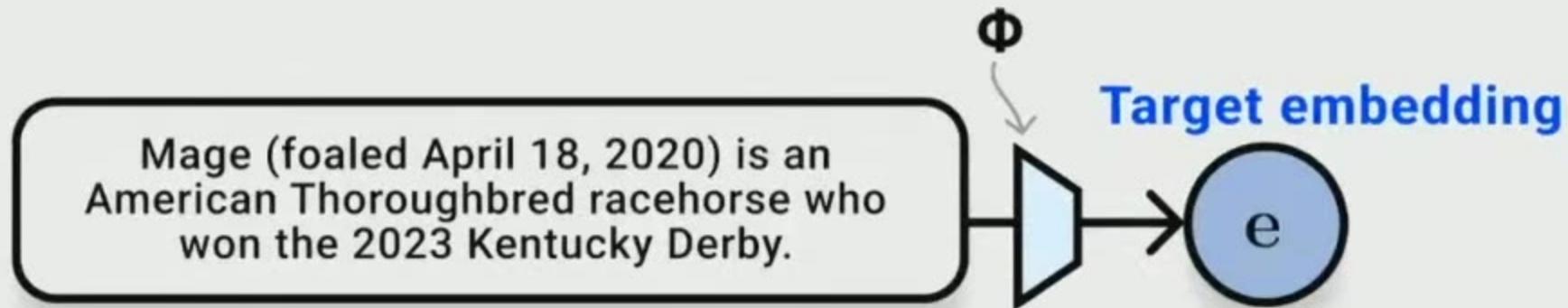


Baseline approach: conditional generation



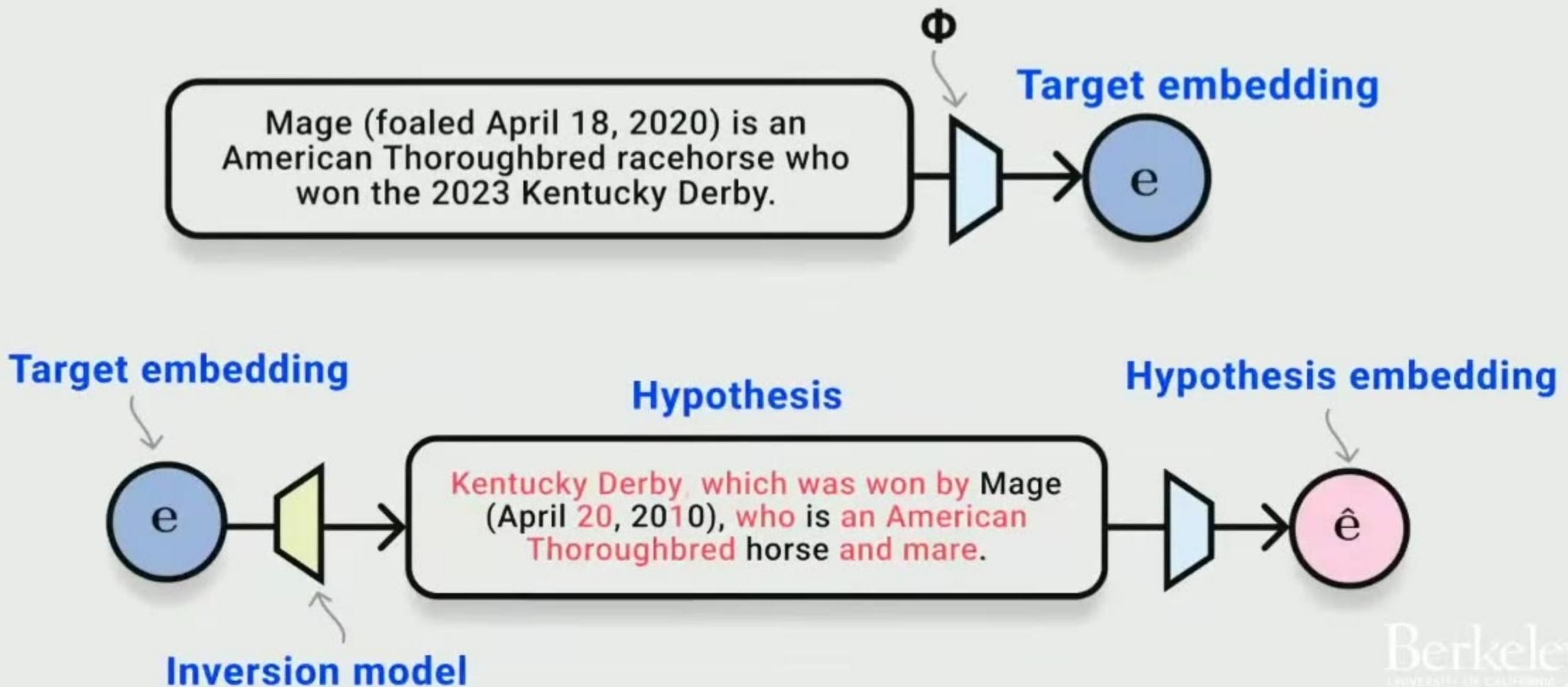


Observation: generated text has a *different* embedding



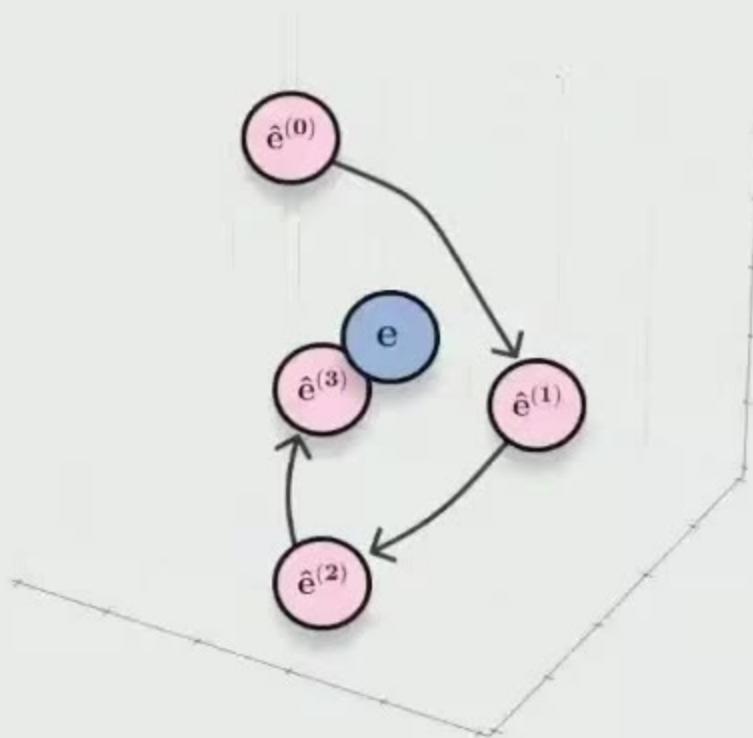


Observation: generated text has a *different* embedding





Idea: iteratively refine and re-embed

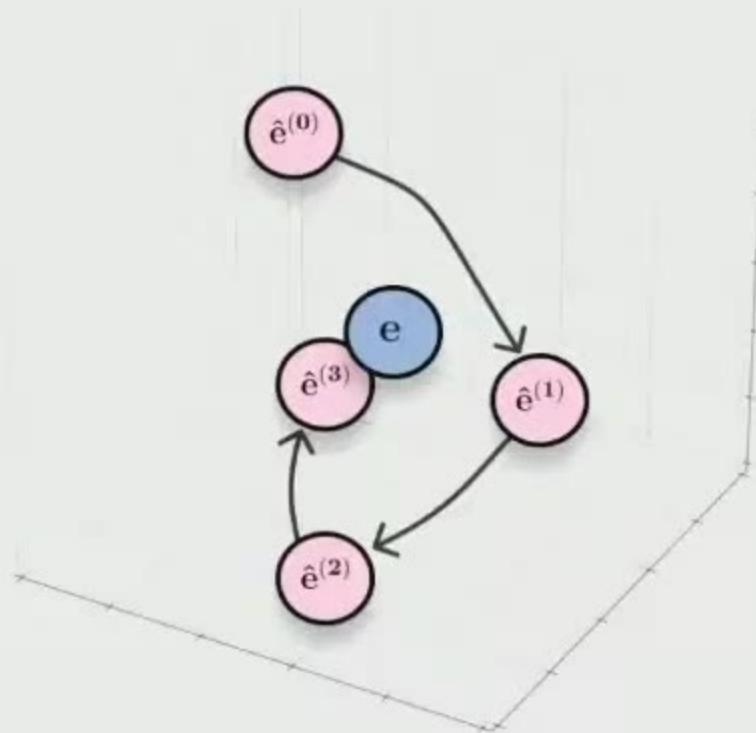


Inspired by approaches to *Iterative Refinement* (Lee et al., 2018; Ghazvininejad et al., 2019; Welleck et al. 2022)



Verizon Event ...

Idea: iteratively refine and re-embed



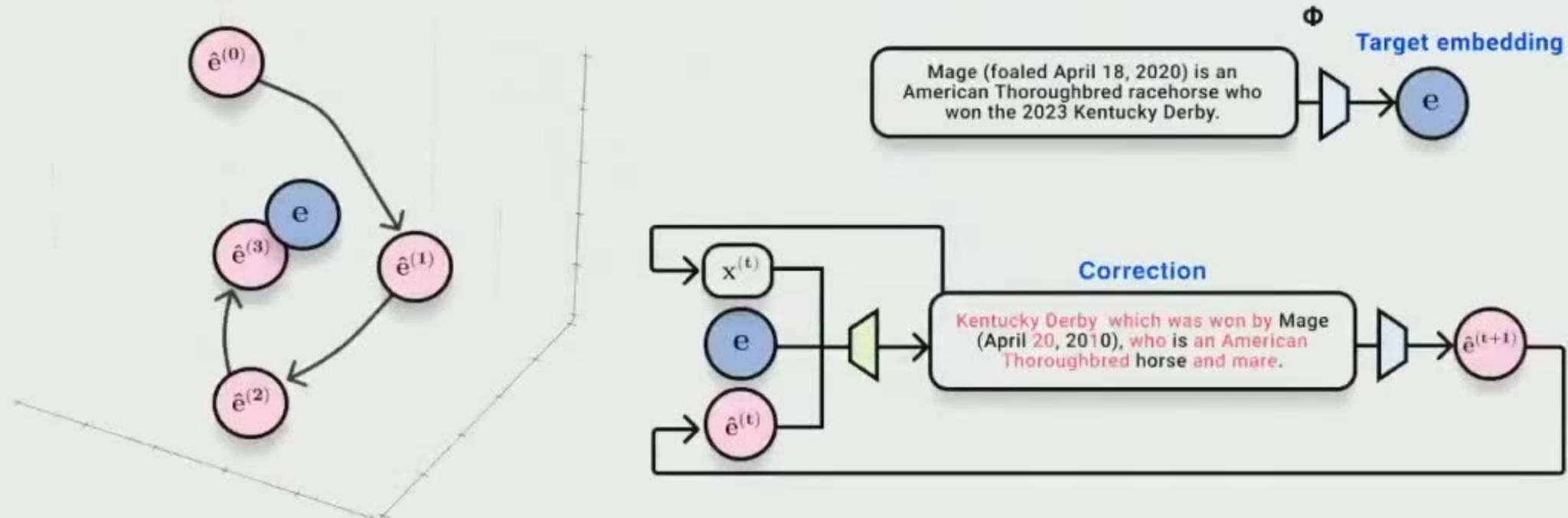
Inspired by approaches to *Iterative Refinement* (Lee et al., 2018; Ghazvininejad et al., 2019; Welleck et al., 2022)

Berkeley
Powered by Zoom



Verizon Event ...

Method: iteratively refine and re-embed (`vec2text`,





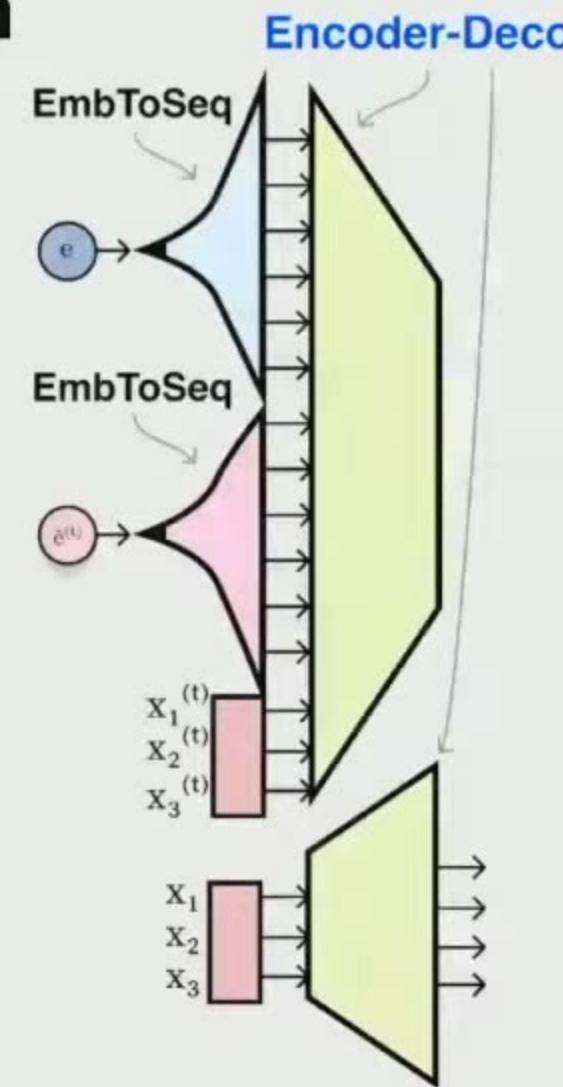
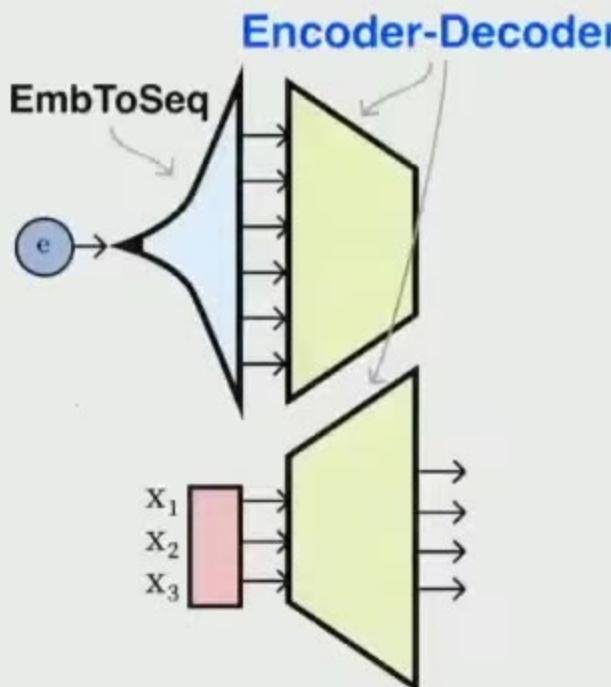
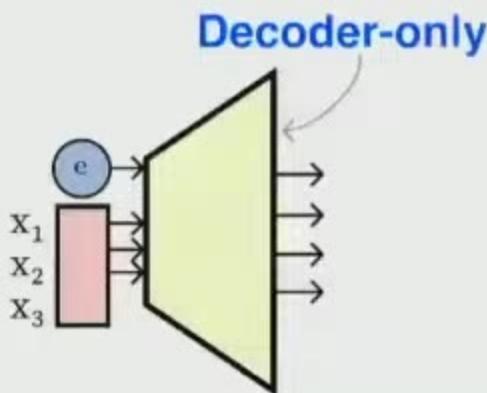
Method: iterative refinement (`vec2text`)

$$\begin{aligned}
 & \text{target} \\
 & \text{embedding} \\
 p(x^{(t+1)} | e) &= \sum_{x^{(t)}} p(x^{(t)} | e) \cdot p(x^{(t+1)} | e, x^{(t)}, \hat{e}^{(t)}) \\
 \hat{e}^{(t)} &= \phi(x^{(t)}) \\
 & \text{embedder} \\
 & \text{hypothesis (t)}
 \end{aligned}$$

The diagram illustrates the iterative refinement process. It starts with a "target embedding" (purple text) which is used to calculate the probability of the hypothesis at time t+1 given the embedding e. This probability is the product of the probability of the hypothesis at time t given e, and the probability of the hypothesis at time t+1 given the hypothesis at time t and the embedding e. The hypothesis at time t is generated by an "embedder" (blue text) from the hypothesis at time t (red text).



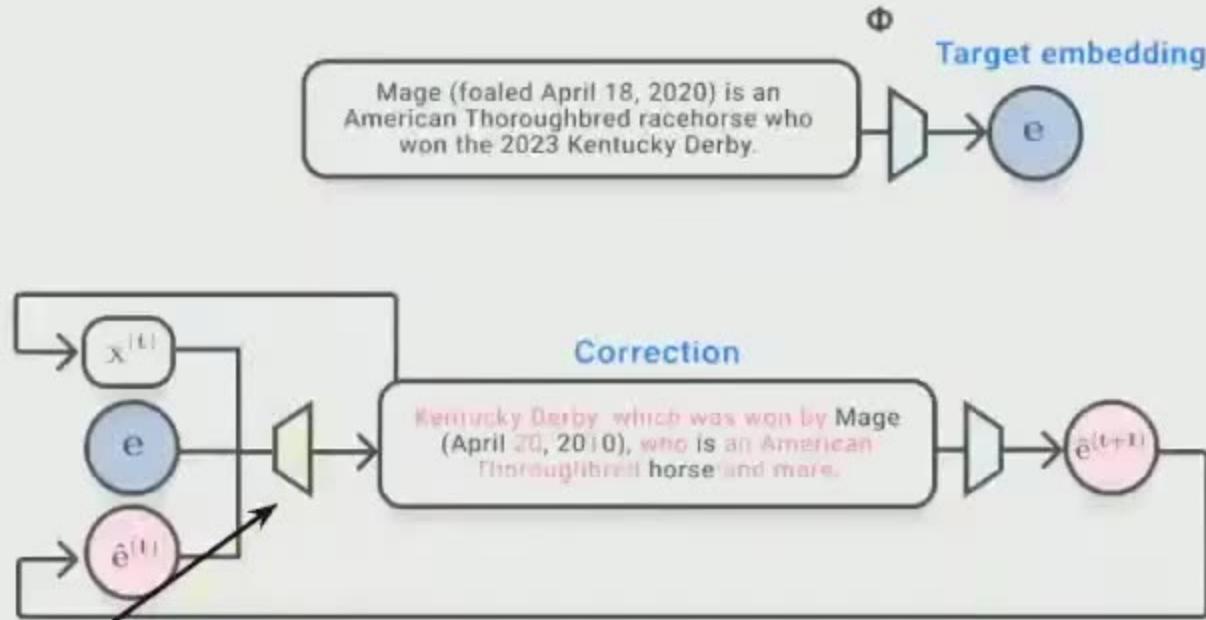
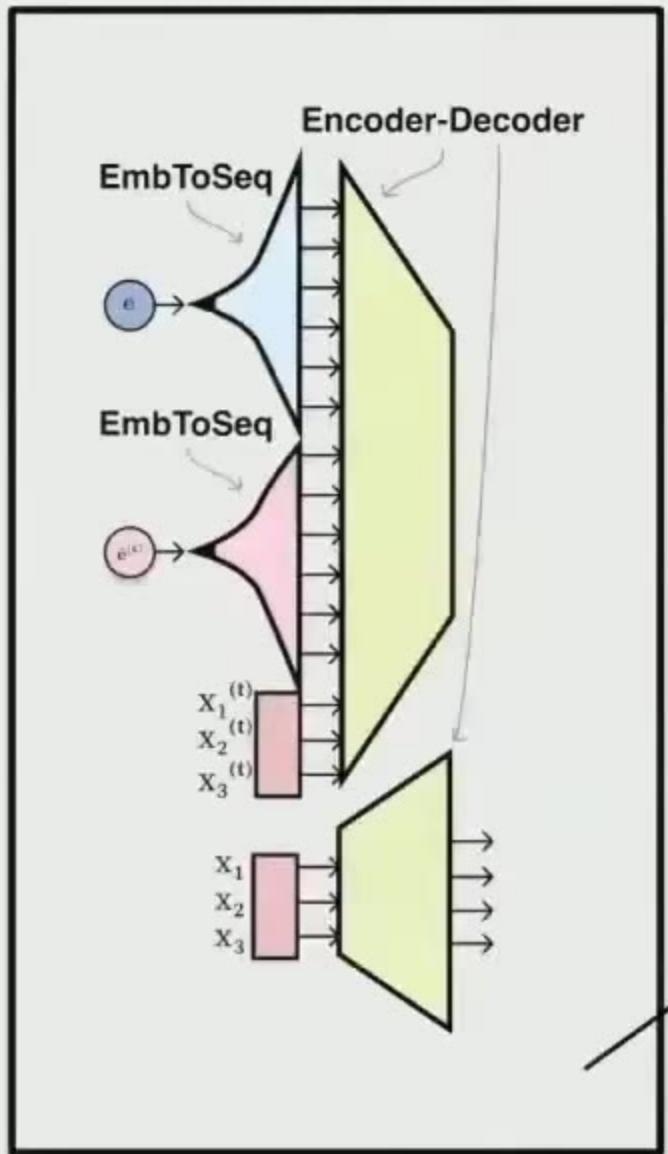
Our approach: iterative generation (architecture)



Prior work

Initial approach

vec2text





Original text

Hypothesis (Round 0)

Embedding





Metrics

- **BLEU score** – approximately how well does the reconstructed text match the original?
- **Exact match** – how often did we reconstruct the original text perfectly?

Baselines

- Bag-of-words model
- Decoder-only architecture
- Correction with fewer steps

Results – Natural Questions

(GTR embeddings || 32 tokens)



[Exact match]

0.0

0.0

[BLEU Score]

0.3

1.0

Bag-of-words
(prior work)

GPT-2 Decoder
(prior work)

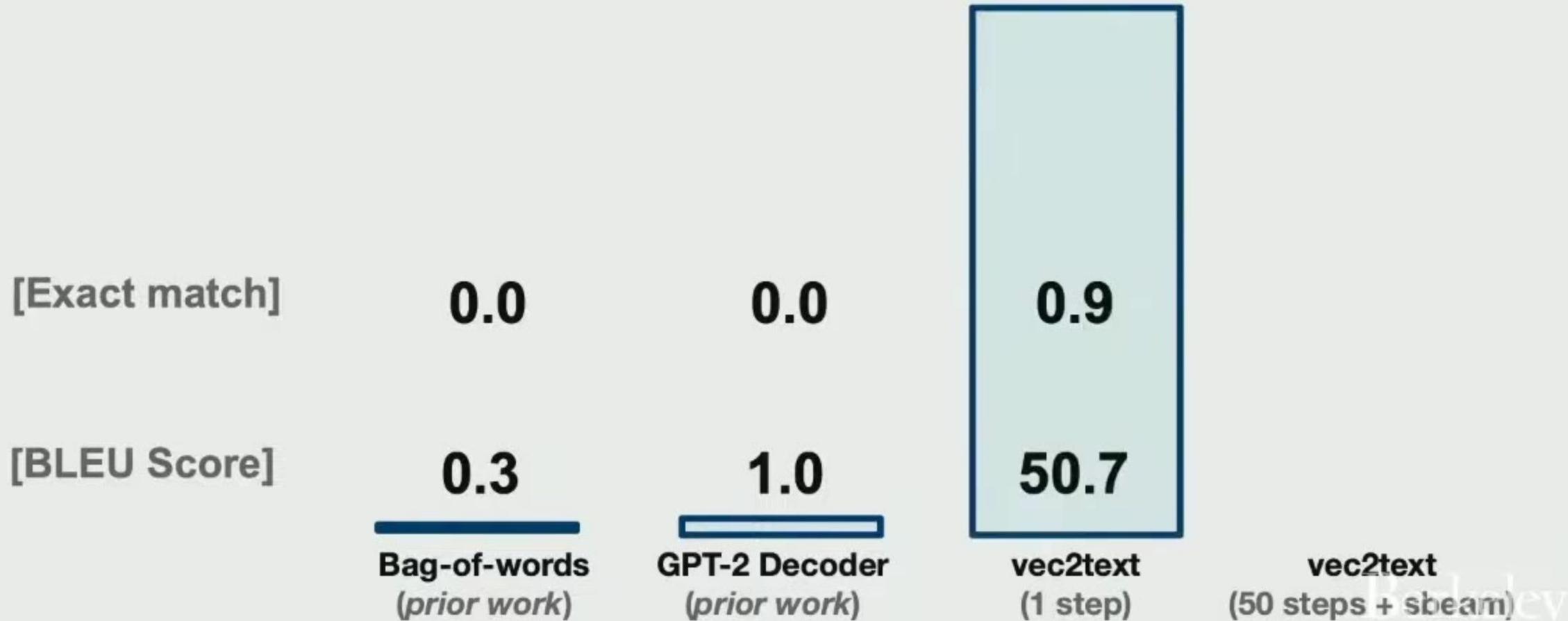
vec2text
(1 step)

vec2text
(50 steps + sbeam)



Results – Natural Questions

(GTR embeddings || 32 tokens)



Results – Natural Questions

(GTR embeddings || 32 tokens)



[Exact match]

0.0

0.0

0.9

92.0

[BLEU Score]

0.3

1.0

50.7

97.1

Bag-of-words
(prior work)

GPT-2 Decoder
(prior work)

vec2text
(1 step)

vec2text
(50 steps + sbeam)
databricks.com



Results – MSMARCO

(OpenAI embeddings || 128 tokens)

[Exact match]

0.6

1.4

3.4

8.0

[BLEU Score]

17.0

29.9

44.4

55.0

Base model
(0 steps)

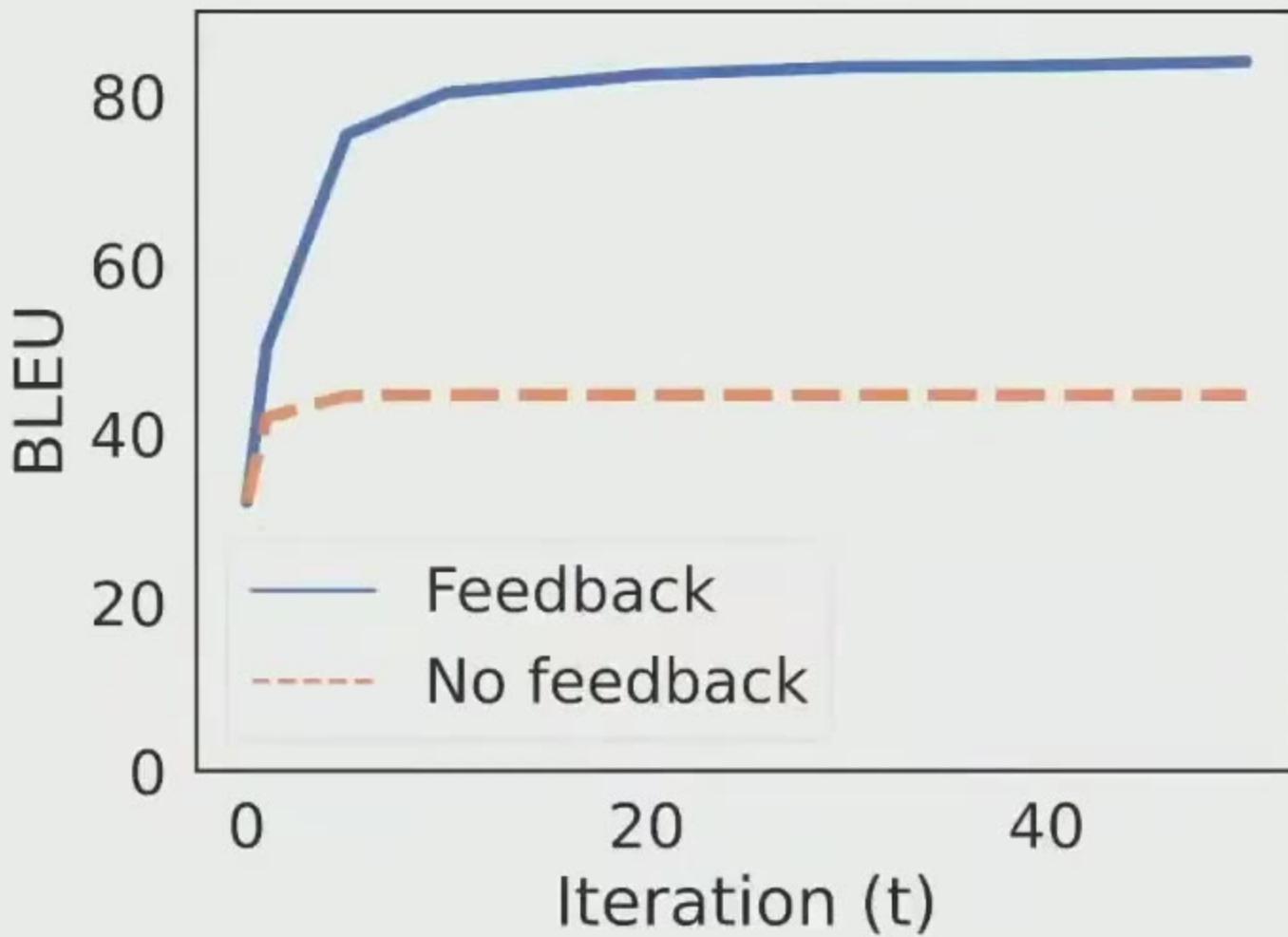
vec2text
(1 step)

vec2text
(50 steps)

vec2text
(50 steps + sbeam)

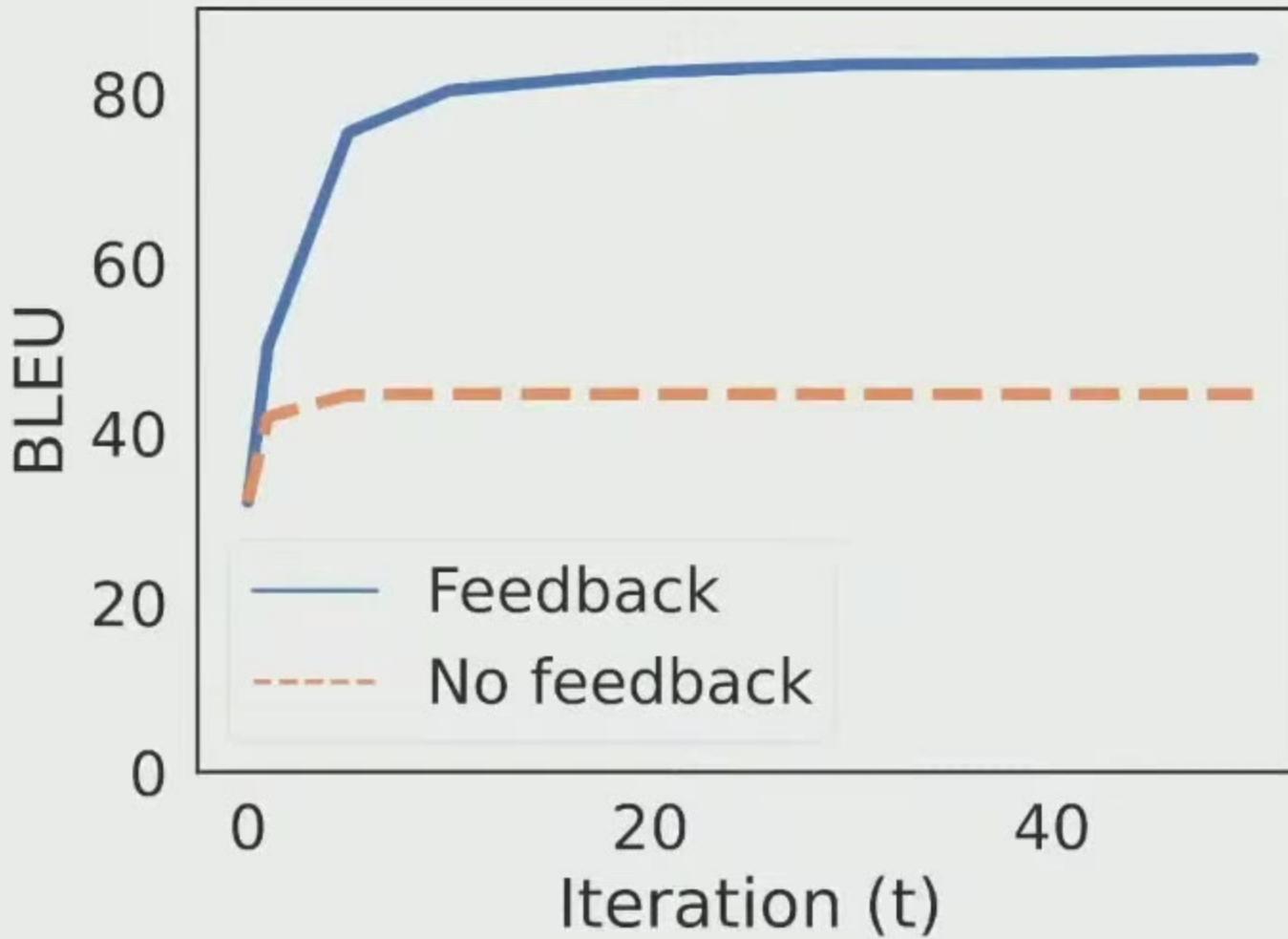


Inversion with and without feedback



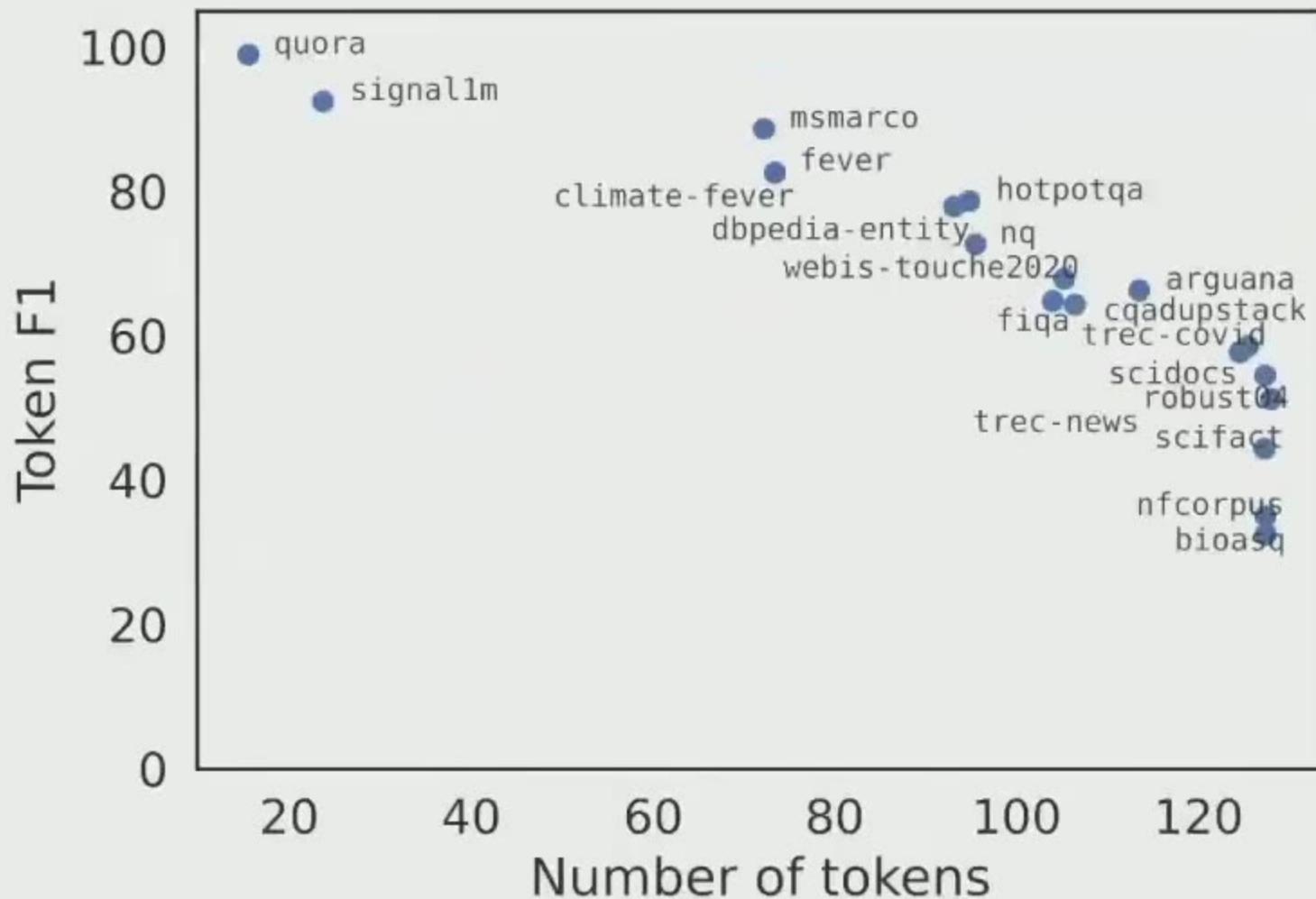


Inversion with and without feedback

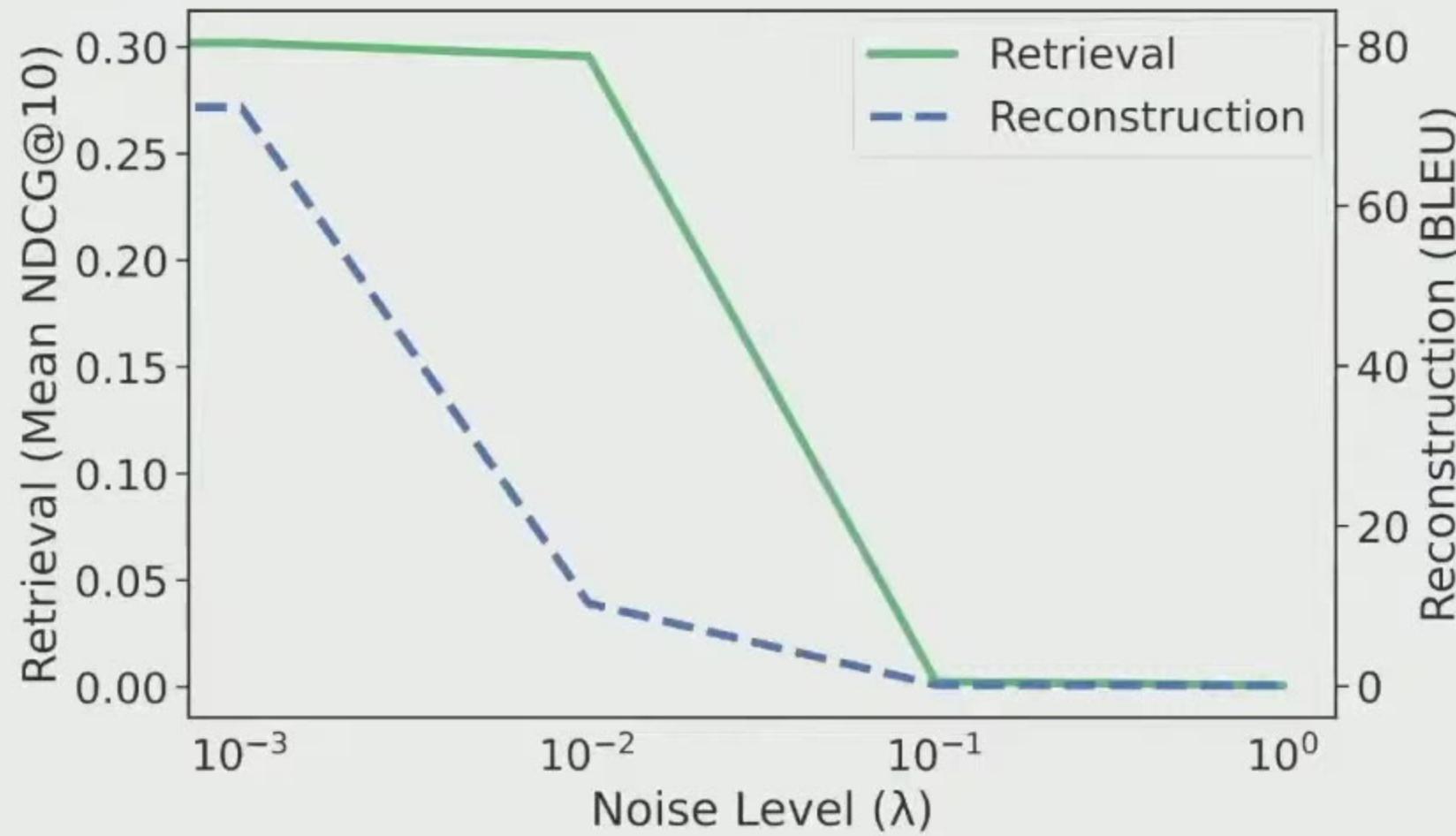




Inversion performance vs. sequence length – BEI



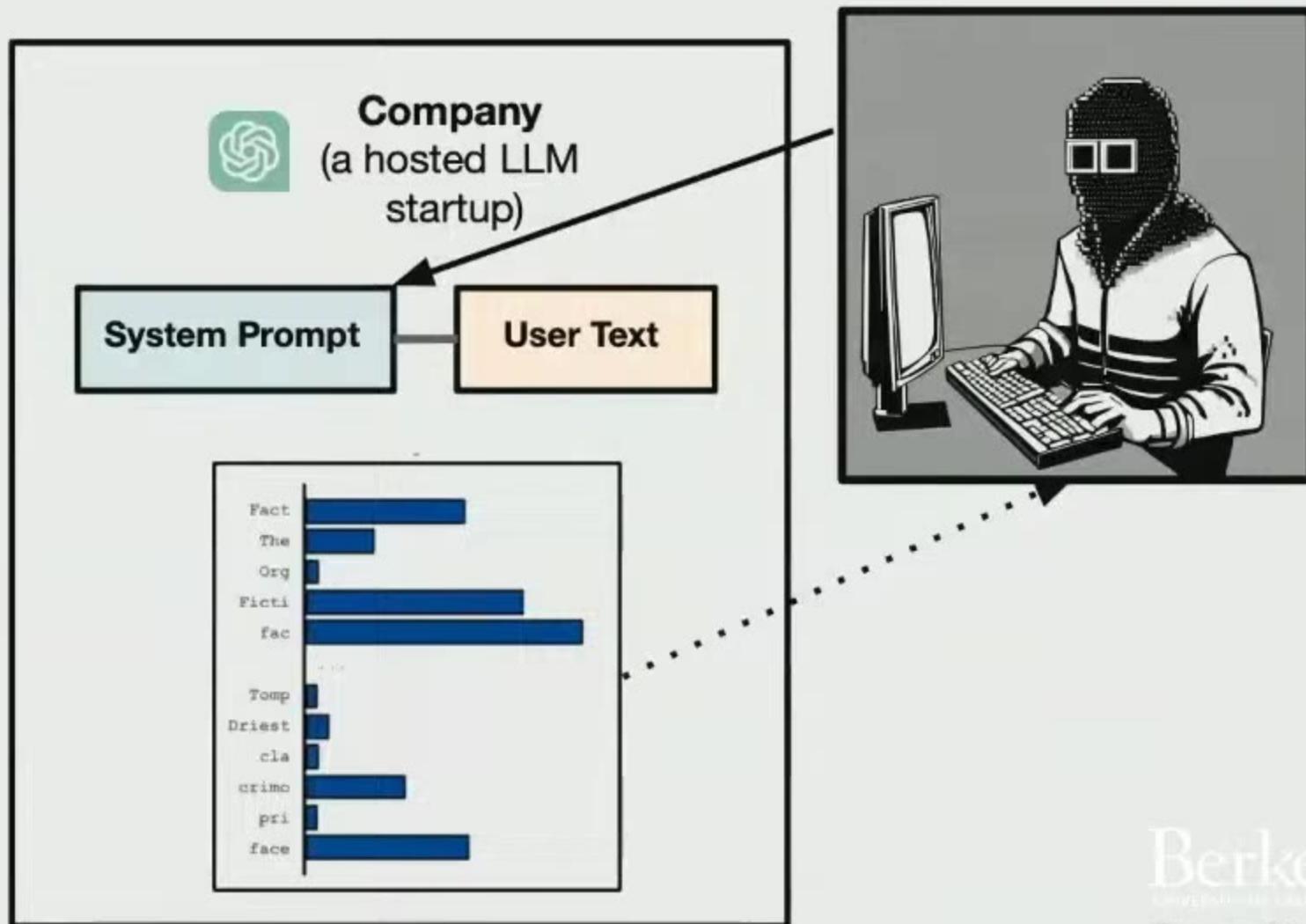
Defending against inversion





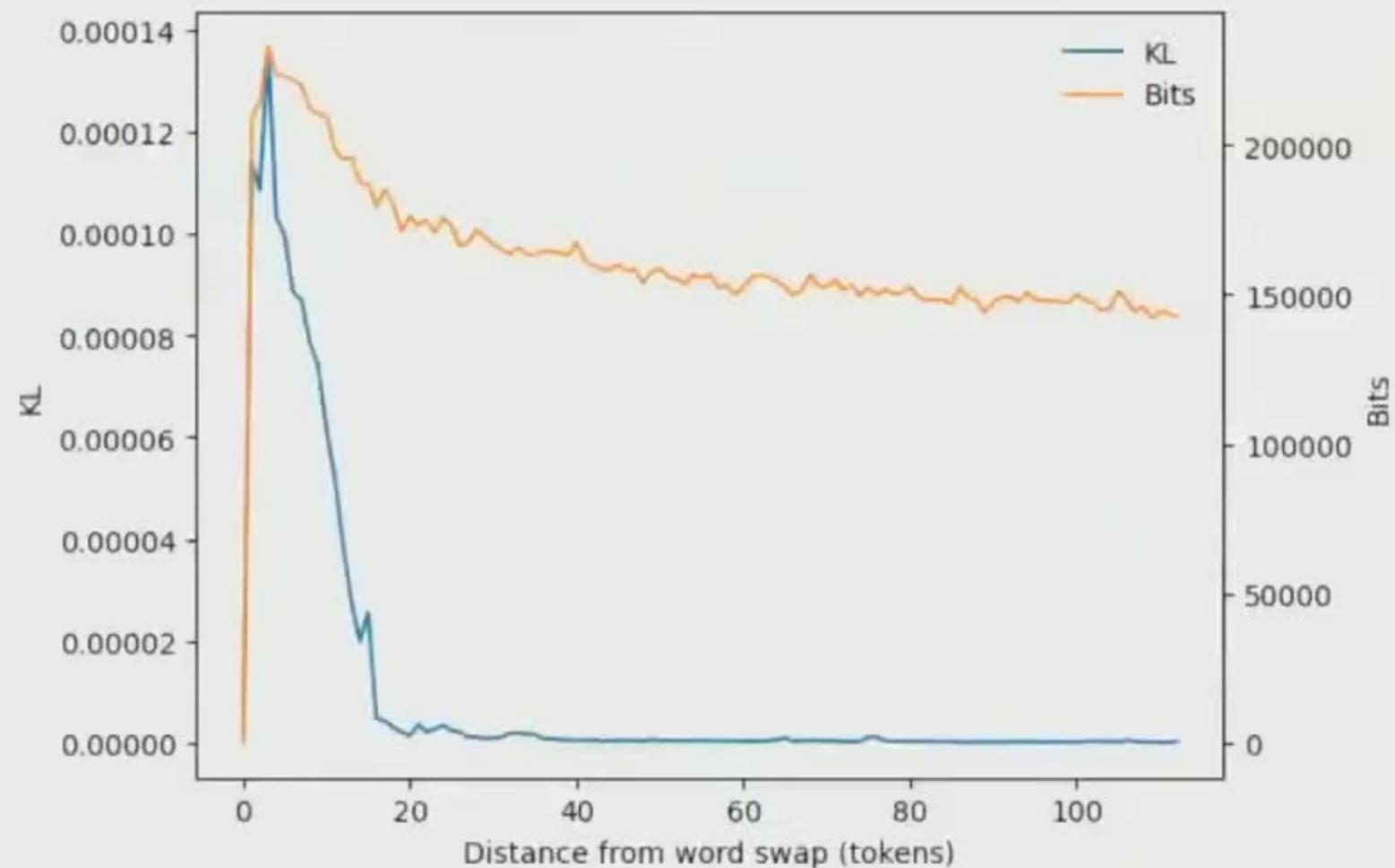
Threat model

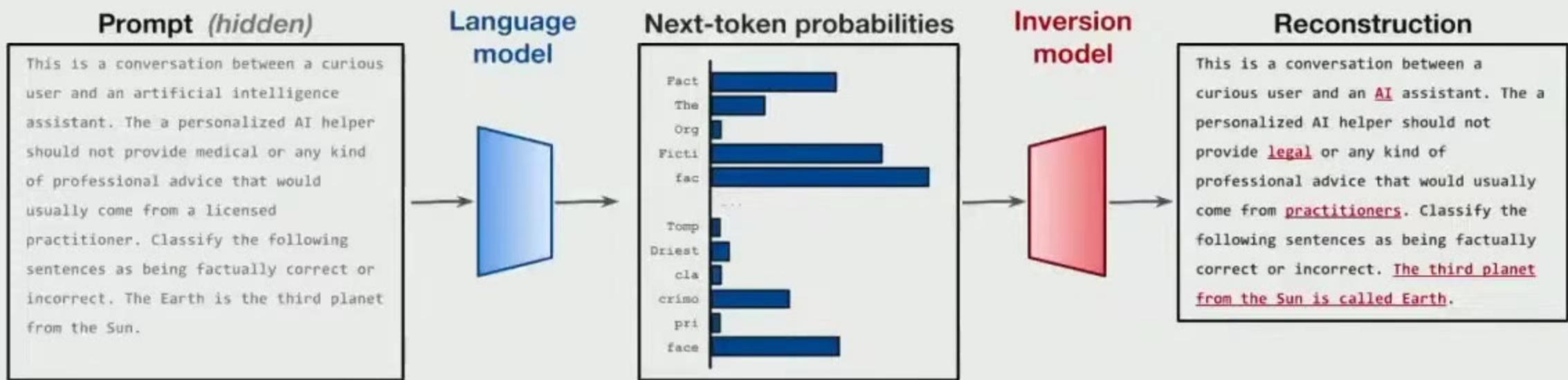
What can a *bad actor* learn from just next word distribution of text?





Distributions are similar. Vectors are different





Results – Instructions 2M

(Llama2 Chat || 64)

Verizon Event ...

[Exact match]

0.0

0.0

0.0

23.4

[BLEU Score]

25.5

6.1

14.9

58.3

Sample
Inverter

Few-shot
(GPT4)

Jailbreak
(oracle)

Distribution
Inversion



Powered by Zoom



LMs Inside Out

- Iterative methods can recover embeddings
- Embeddings are not compressive

Language models

- LM distributions are information dense
- Jailbreak protected models can still be recovered



[jxmorris12/vec2text](https://github.com/jxmorris12/vec2text)



Berkeley
UNIVERSITY OF CALIFORNIA
Powered by Zoom



LMs Inside Out

- Iterative methods can recover embeddings
- Embeddings are not compressive

Language models

- LM distributions are information dense
- Jailbreak protected models can still be recovered



[jxmorris12/vec2text](https://github.com/jxmorris12/vec2text)



Berkeley
UNIVERSITY OF CALIFORNIA

Powered by Zoom



LMs Inside Out

- Iterative methods can recover embeddings
- Embeddings are not compressive

Language models

- LM distributions are information dense
- Jailbreak protected models can still be recovered



[jxmorris12/vec2text](https://github.com/jxmorris12/vec2text)



Berkeley
UNIVERSITY OF CALIFORNIA

Powered by Zoom

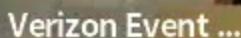


Berkeley Center for Responsible
Decentralized Intelligence

Summit on Responsible Decentralized Intelligence — Future of Decentralization and AI

August 6, 2024
Verizon Center, NYC

Berkeley
UNIVERSITY OF CALIFORNIA
Powered by Zoom



Verizon Event ...

Berkeley Center for Responsible
Decentralized Intelligence

Summit on Responsible Decentralized Intelligence — Future of Decentralization and AI

August 6, 2024
Verizon Center, NYC

Berkeley
UNIVERSITY OF CALIFORNIA
Powered by Zoom

Berkeley



Verizon Event ...

Session I: Open Source AI
Berkeley

Center for Responsible,
Decentralized Intelligence

Building a Global AI Safety Benchmark

Summit on Responsible Decentralized Intelligence
— Future of Decentralization and AI

Mala Kumar

AI Safety

MLCommons

August 6, 2024

Verizon Center, NYC

Berkeley
UNIVERSITY OF CALIFORNIA
Powered by Zoom



Session I: Open Source AI

Building a Global AI Safety Benchmark

Mala Kumar

Director of Program Management
AI Safety
MLCommons





Session I: Open Source AI

Building a Global AI Safety Benchmark

Mala Kumar

Director of Program Management
AI Safety
MLCommons

