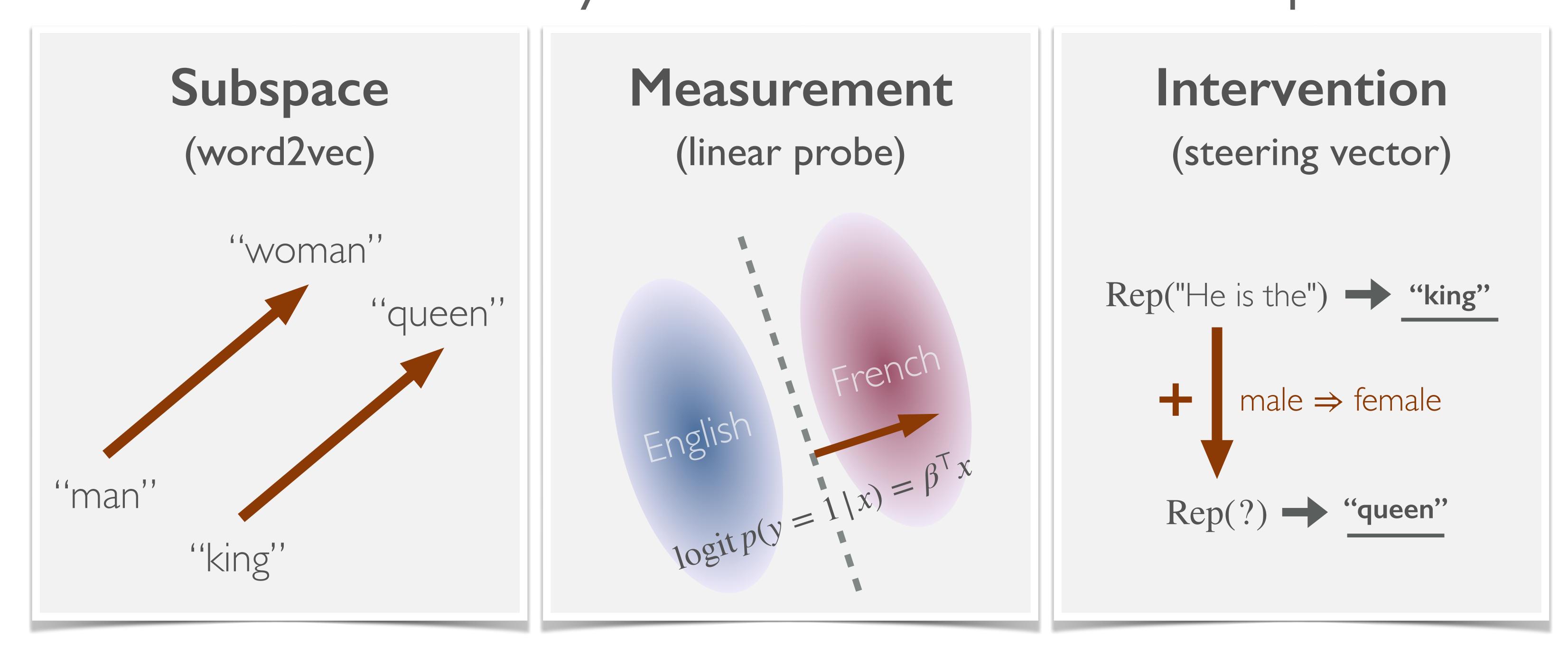
## Linear Representation Hypothesis & Geometry of LLMs

Kiho Park<sup>1</sup>, Yo Joong (YJ) Choe<sup>2</sup>, Victor Veitch<sup>1,2</sup>

<sup>1</sup>Department of Statistics, <sup>2</sup>Data Science Institute, University of Chicago

#### What Does "Linear" Even Mean?

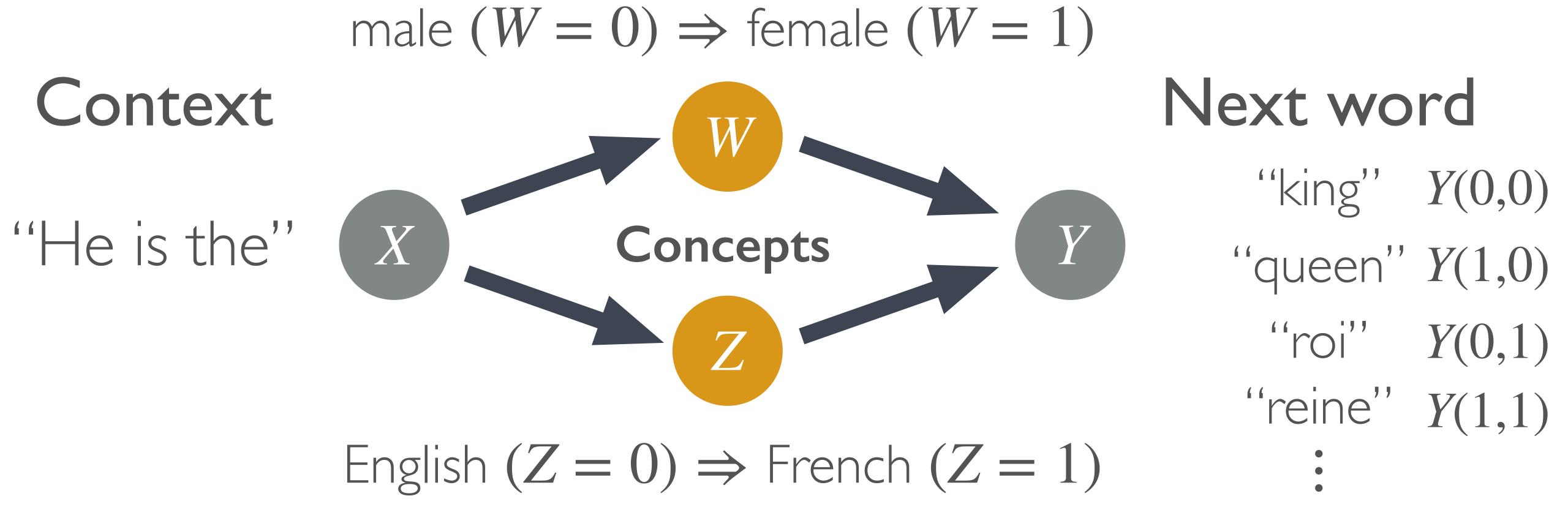
There are three natural ways to define the notion of linear representation.



Problem: It is not clear how these ideas relate to each other, nor which is the right notion of linear representation.

## Concepts in LLMs

A concept W is defined by counterfactual outputs Y(W=0), Y(W=1). Concepts W and Z are causally separable if Y(w,z) is well-defined.



LLMs generate the next word using the softmax distribution.

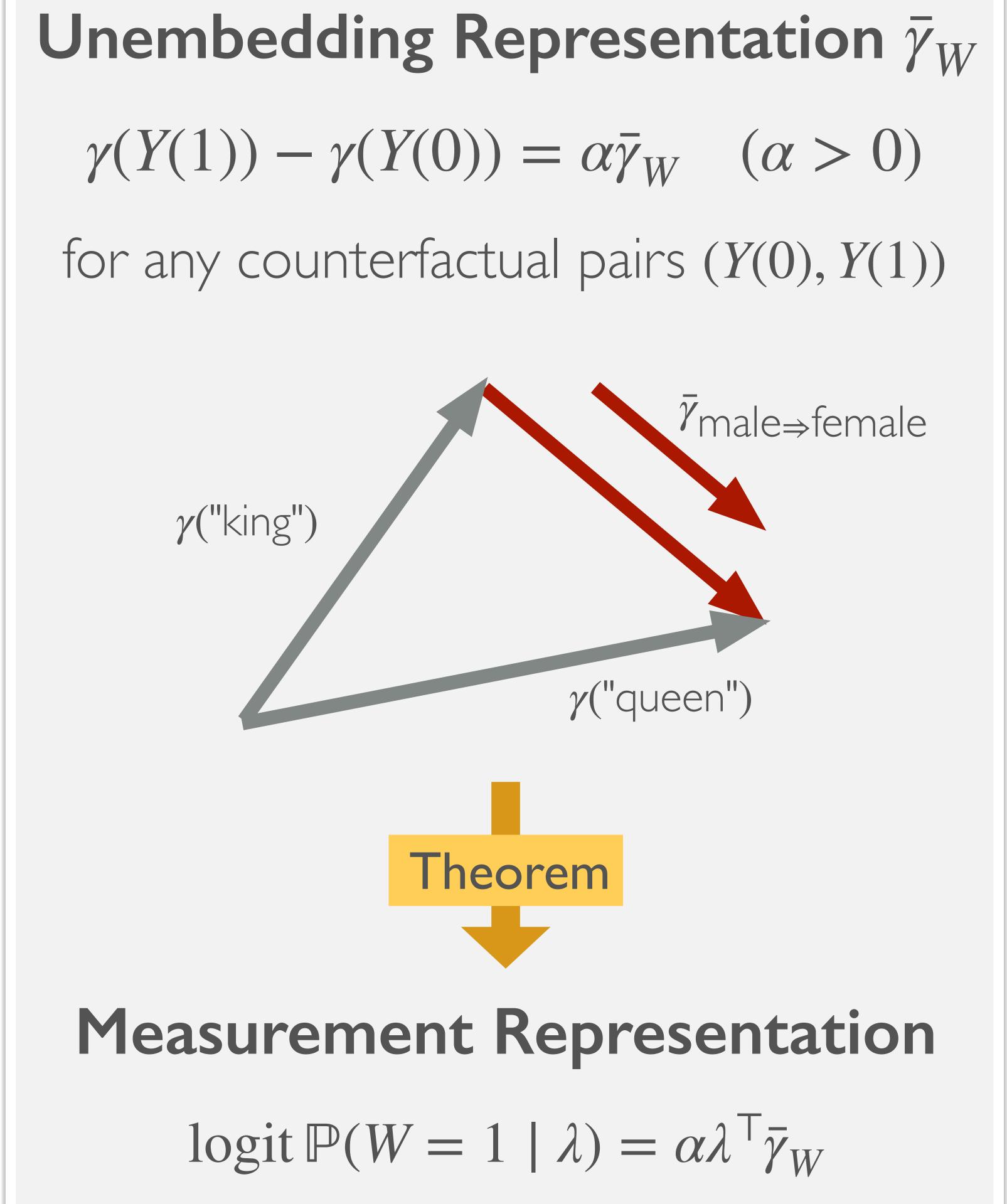
Embedding  $\lambda(x) \in \mathbb{R}^d$ 

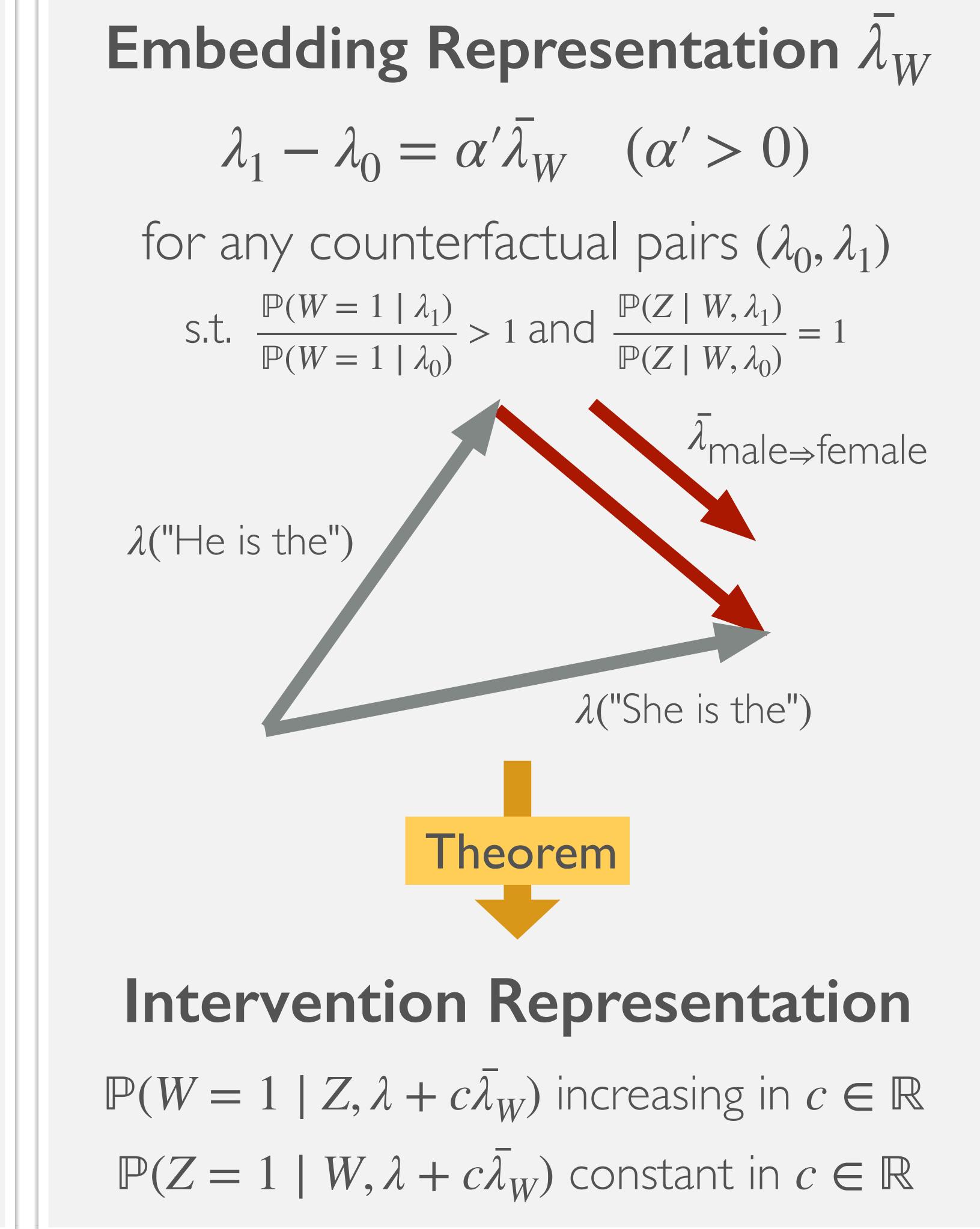
Softmax  $\mathbb{P}(y \mid x) \propto \exp(\lambda(x)^{\mathsf{T}} \gamma(y))$ 

Unembedding  $\gamma(y) \in \mathbb{R}^d$ 

## Formalizing Linear Representation Hypothesis

We first formalize the subspace notions of linear representations, then use softmax structure to connect them to measurement and intervention.





#### Problems:

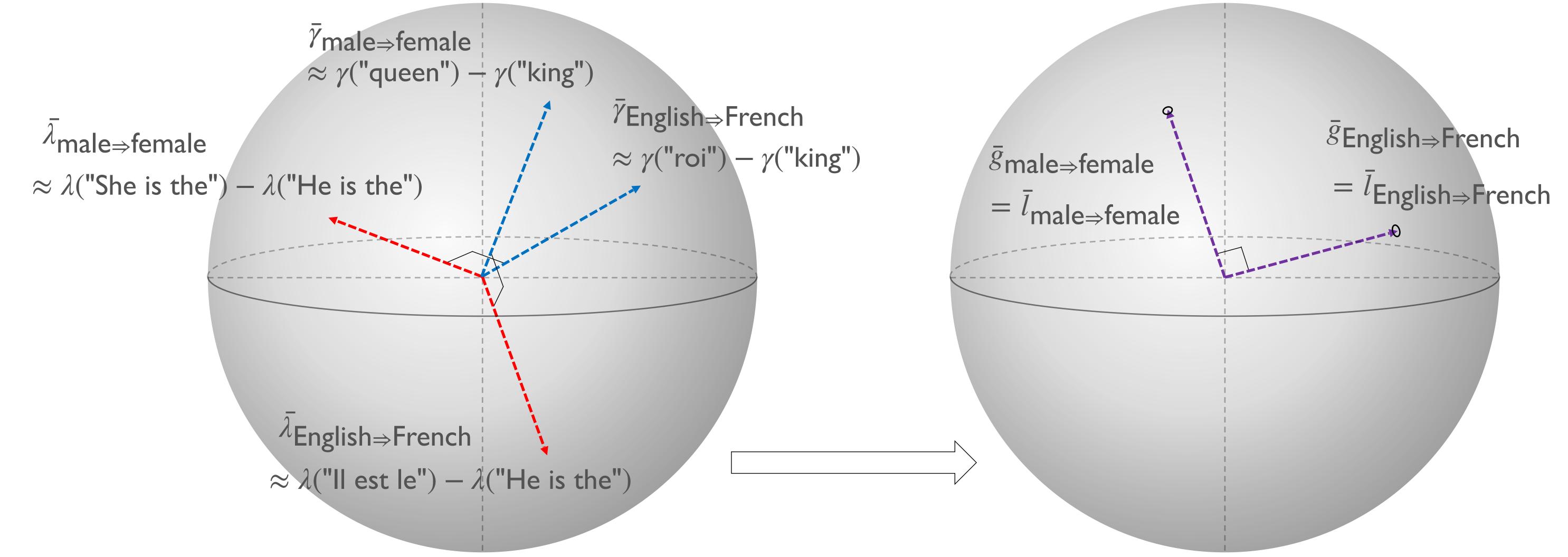
How do the unembedding and embedding representations relate? What is the right inner product for the representation space?

#### Causal Inner Product

Definition: Whenever W and Z are causally separable,  $\langle \bar{\gamma}_W, \bar{\gamma}_Z \rangle_C = 0$ .

Theorem

This unifies the unembedding and embedding representations via  $\langle \bar{\gamma}_W, \cdot \rangle_C = (\bar{\lambda}_W)^{\top}$ . (This is the Riesz isomorphism.)



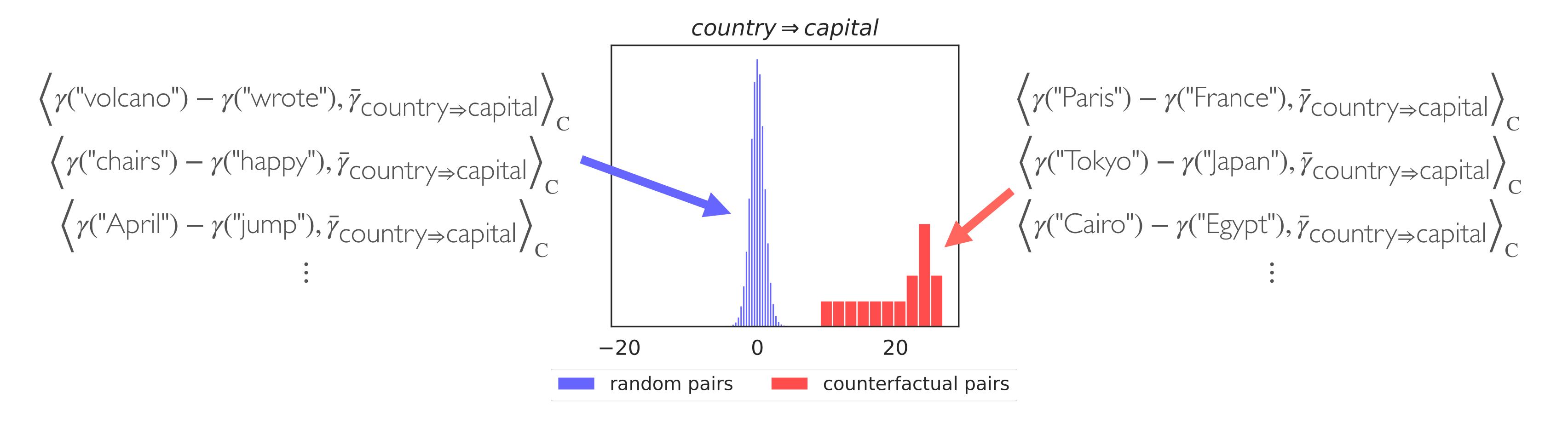
We can estimate the causal inner product as  $\langle \bar{\gamma}, \bar{\gamma}' \rangle_C = \bar{\gamma}^T \text{Cov}(\gamma)^{-1} \bar{\gamma}'$ .

## Experiments with LLaMA-2 Representations

#### Linear Representations Exist

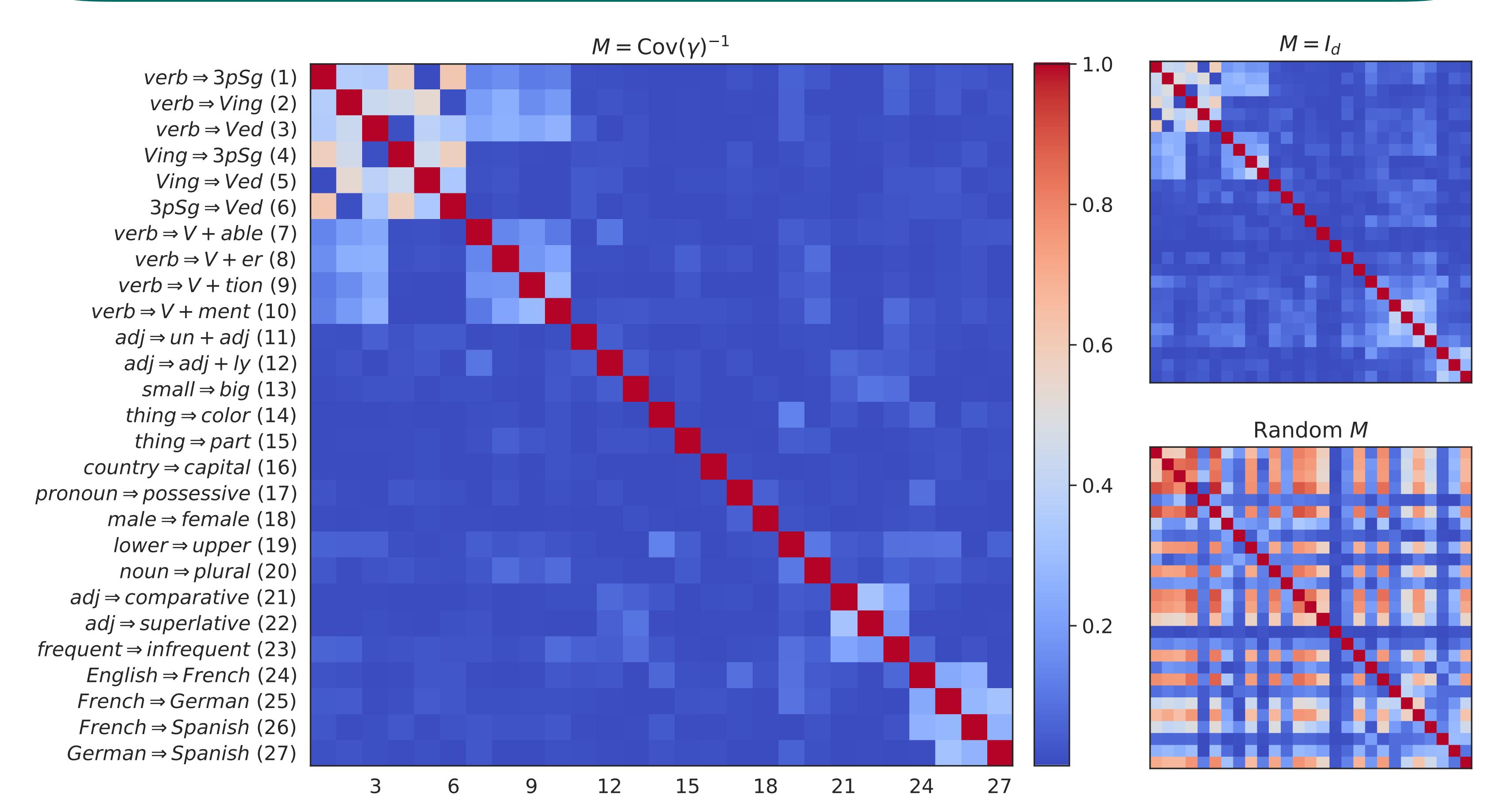
We estimate the unembedding representations for various concepts by using the counterfactual pairs from a word analogy dataset.

$$\bar{\gamma}_W := \frac{\tilde{\gamma}_W}{\langle \tilde{\gamma}_W, \tilde{\gamma}_W \rangle_{\mathrm{C}}} \text{ where } \tilde{\gamma}_W = \frac{1}{n_W} \sum_{i=1}^{n_W} \gamma(y_i(1)) - \gamma(y_i(0))$$



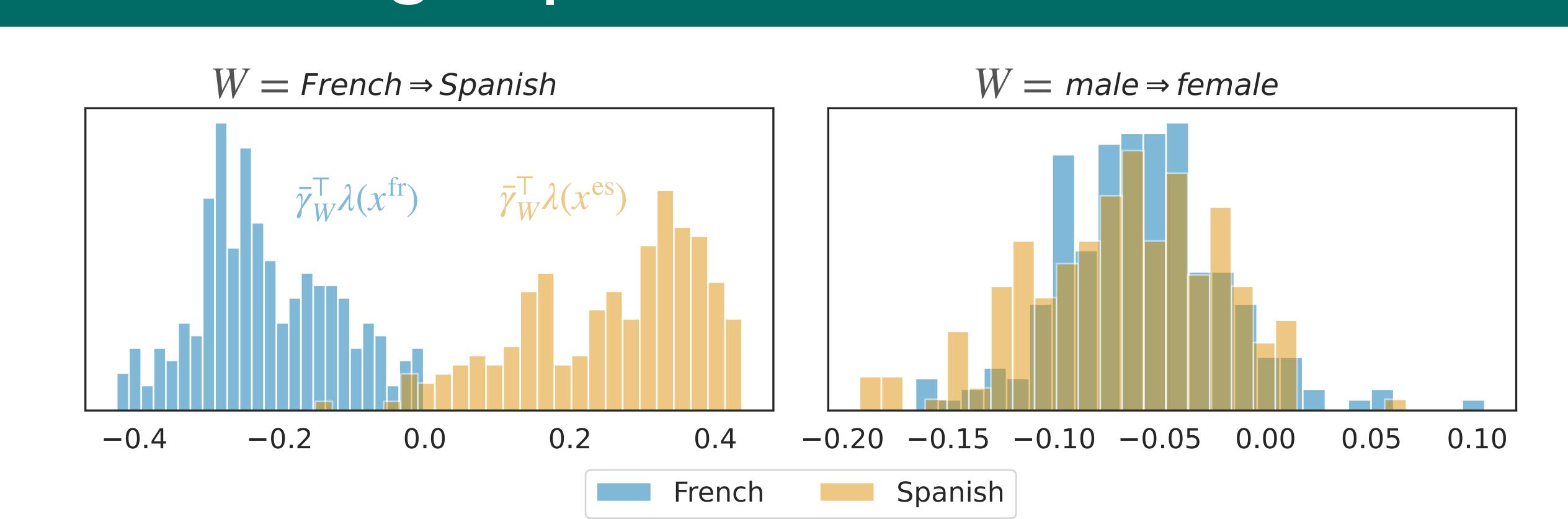
Differences between counterfactual pairs are more parallel to  $\bar{\gamma}_W$  than those between random pairs, supporting the linear representation hypothesis.

# Causally Separable Concepts Are Represented Orthogonally Under the Causal Inner Product



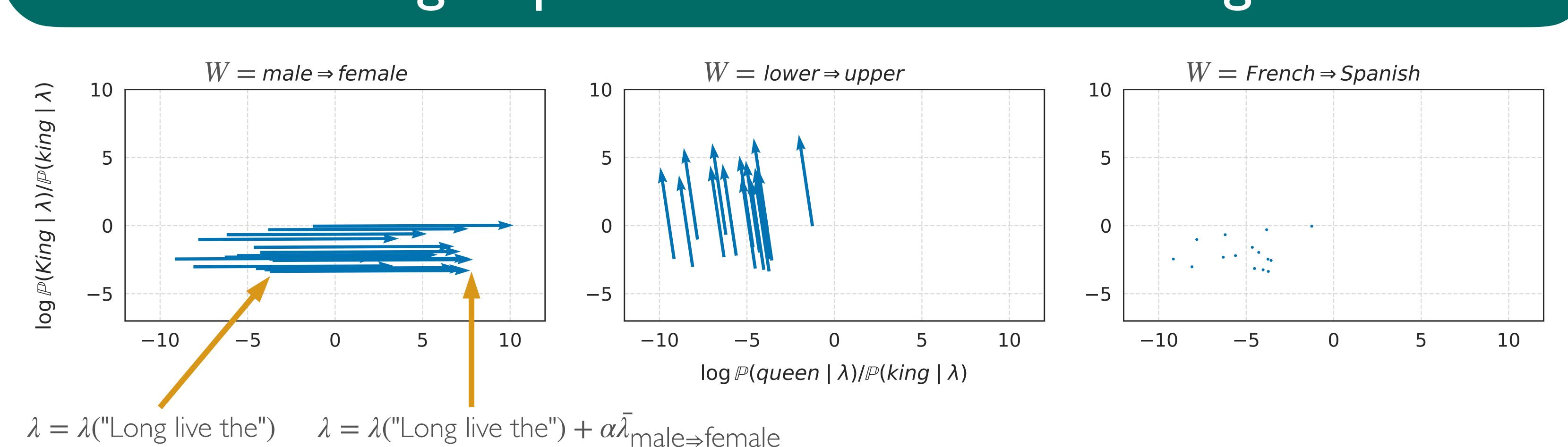
Heatmaps of  $|\bar{\gamma}_W^\top M \bar{\gamma}_W|$  show that the causal inner product between the unembedding representations of causally separable concepts is close to zero. It improves on the naive Euclidean inner product.

### Unembedding Representation Yields Linear Probe



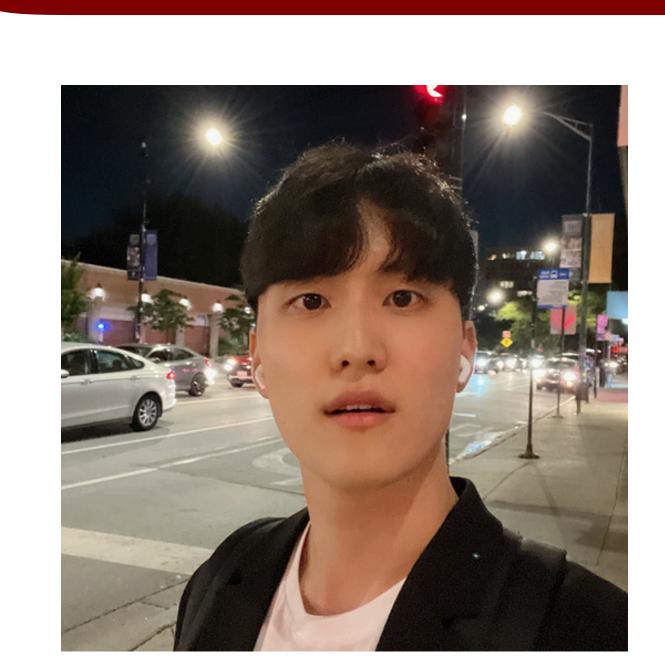
 $\bar{\gamma}_{\text{French} \Rightarrow \text{Spanish}}$  separates the embeddings of French and Spanish contexts, while  $\bar{\gamma}_{\text{male} \Rightarrow \text{female}}$  does not.

#### Embedding Representation Yields Steering Vector



Adding the embedding representation  $\bar{\lambda}_W := \mathrm{Cov}^{-1}(\gamma)\bar{\gamma}_W$  to context embeddings changes the target concept, without changing other causally separable concepts.

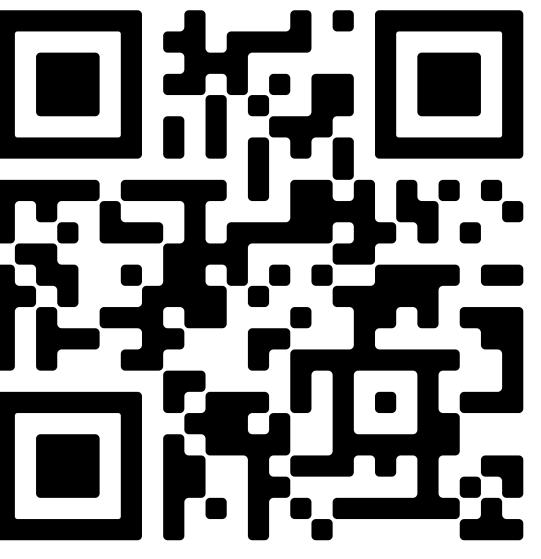
## Find us online! Email: parkkiho@uchicago.edu



Kiho Park







YJ Choe Victor Veitch arXiv:2311.03658