



IMVC 2024

# Kiki or Bouba? Sound Symbolism in Vision-and-Language Models

Morris Alper and Hadar Averbuch-Elor

NeurIPS 2023 ✨ SPOTLIGHT ✨



Next you will be shown two images that were generated by Stable Diffusion using the prompt:

“A 3D rendering of a \_\_\_\_\_-shaped object”

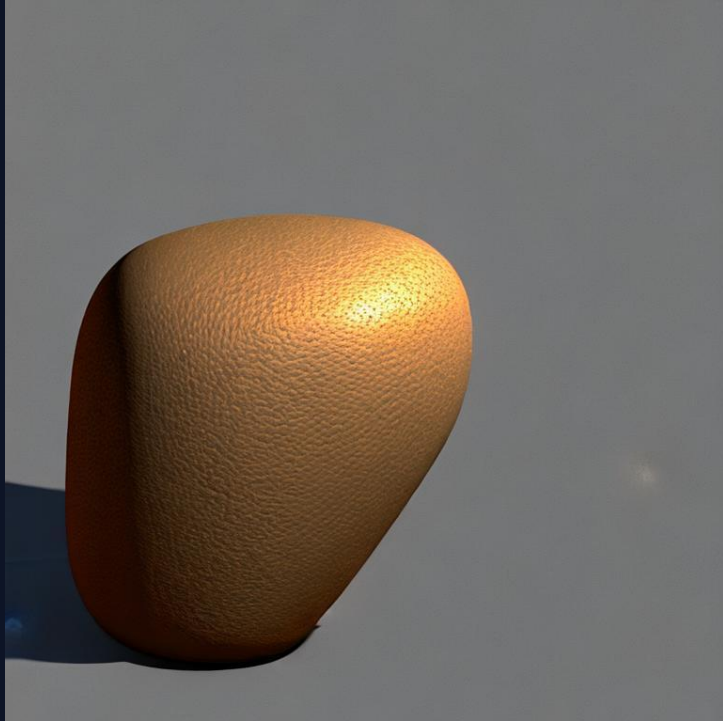
For one the missing word is *bouba* and for the other *kiki*.

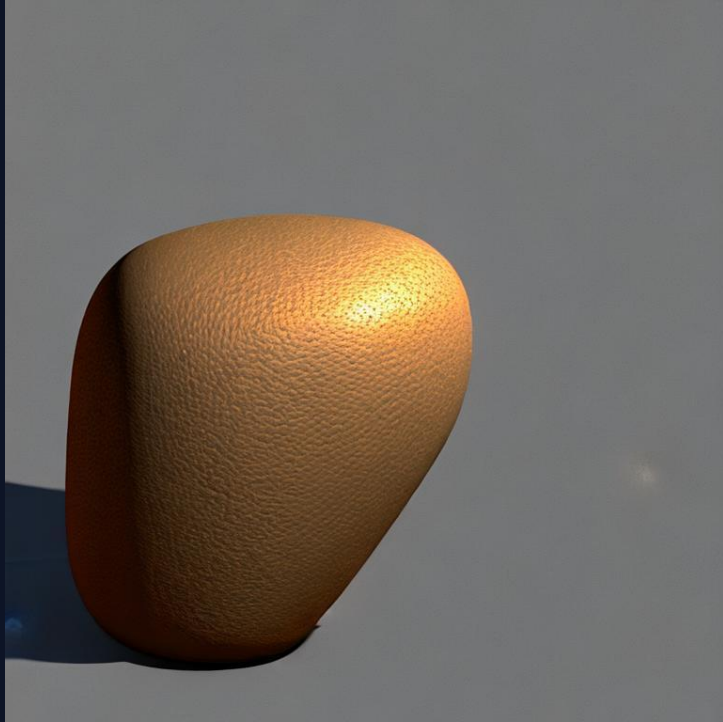
Which is which?



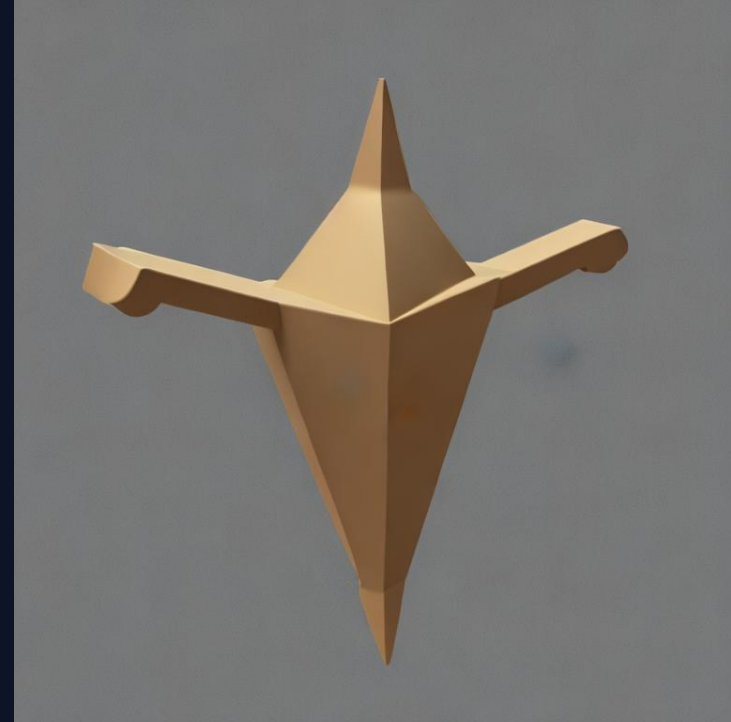
SAGIVTECH

IMVC 2024





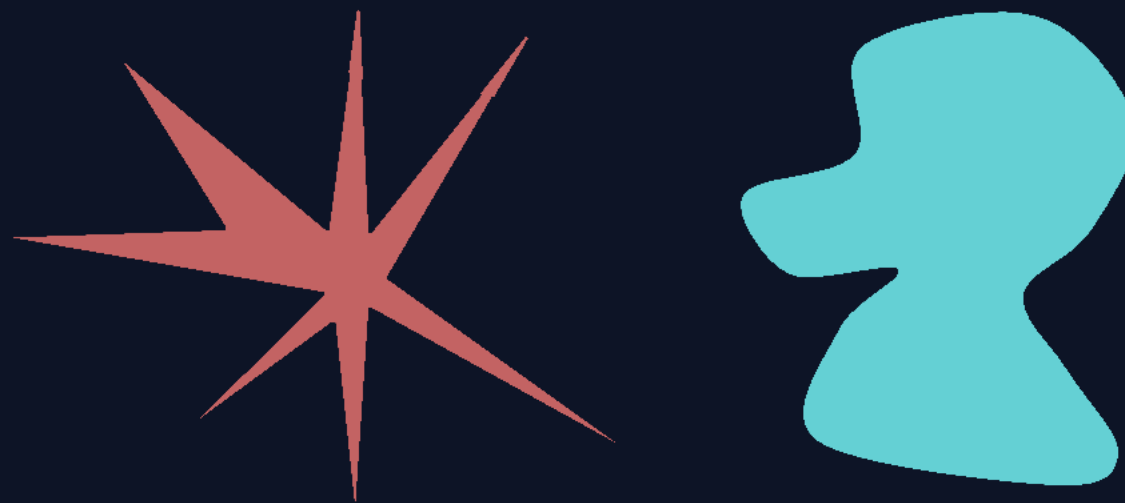
**bouba**



**kiki**

# The Kiki-Bouba Effect

A well-known study in psycholinguistics



Holds across cultures, languages, and even among infants.

# Sound Symbolism

Connections between sound and meaning in language

kapow, glub-glub, boom

crunch/crush/crash/crinkle/crack

mama

אמא (ima)

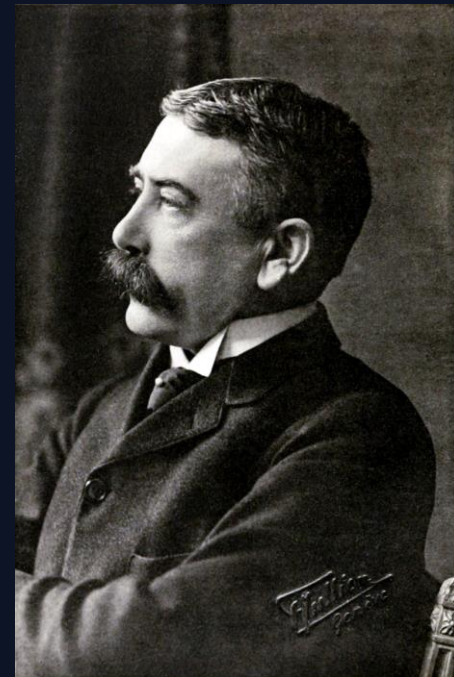
妈妈 (māma)

# Sound Symbolism

Connections between sound and meaning in language

*“Le signe est arbitraire” !*  
(the sign is arbitrary)

**Ferdinand de Saussure**  
(1857–1913)





SAGIVTECH

IMVC 2024

# Sound Symbolism

Connections between sound and meaning in language

***Have VLMs learned sound symbolism?***





SAGIVTECH

IMVC 2024

# Why is this important?

Interpreting the **black box** of foundation VLMs

# Why is this important?

## Interpreting the black box of foundation VLMs



Tewel et al. (CVPR 2022)



Alper et al. (CVPR 2023)

# Why is this important?

Interpreting the black box of foundation VLMs



Tewel et al. (CVPR 2022)



Alper et al. (CVPR 2023)

**New tool for studying language**

Insight into how sound symbolism could be learned



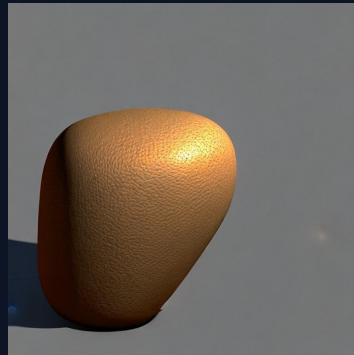
IMVC 2024

# Probing VLMs for Sound Symbolism

# Probing VLMs for Sound Symbolism

Measuring geometric attributes of pseudowords

rounder



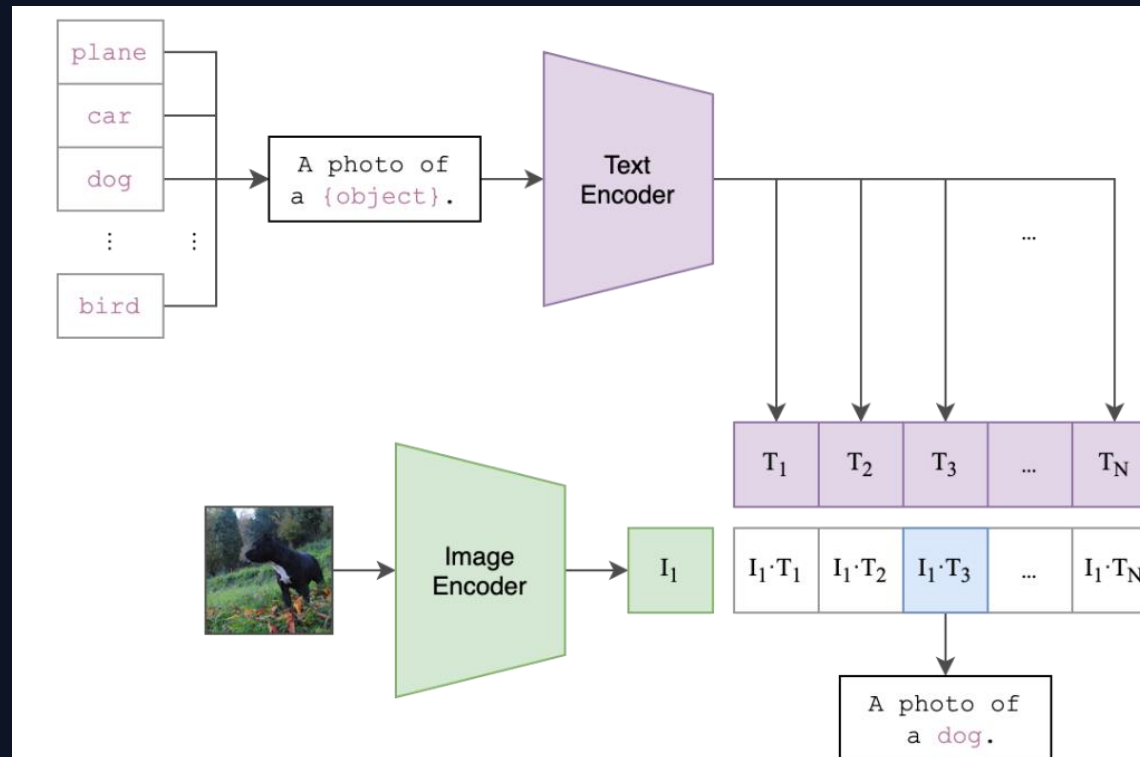
**bouba**

sharper



**kiki**

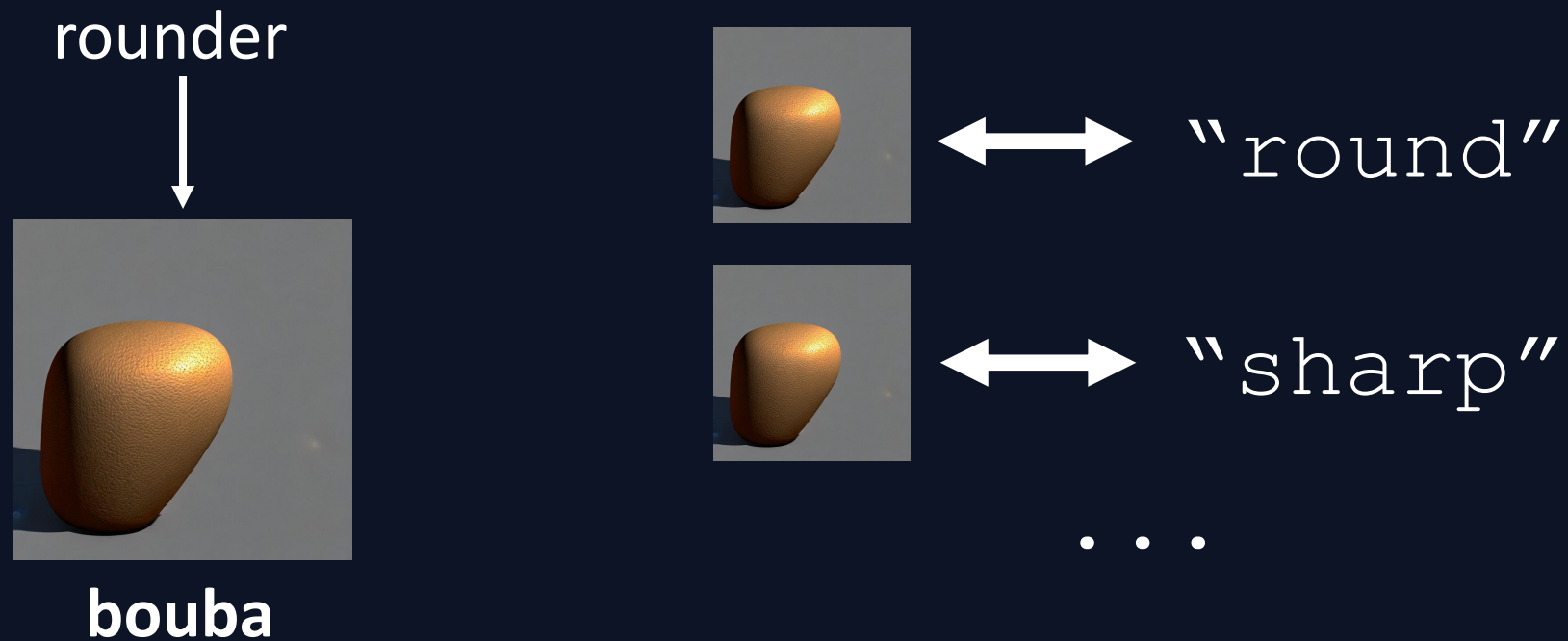
# CLIP: Embeds images and text in a shared semantic space

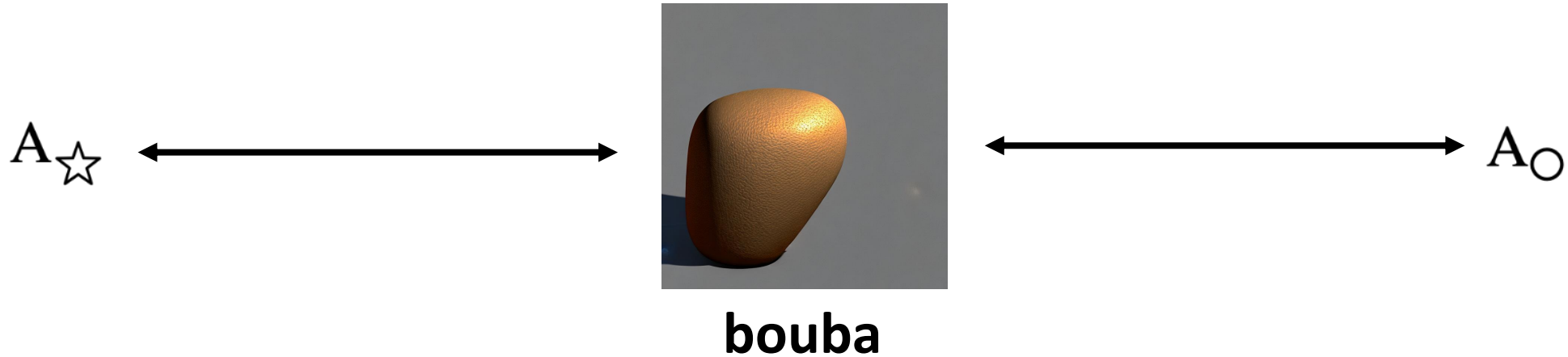


Radford et al. (ICML 2021)

# Probing VLMs for Sound Symbolism

Measuring geometric attributes of pseudowords with CLIP





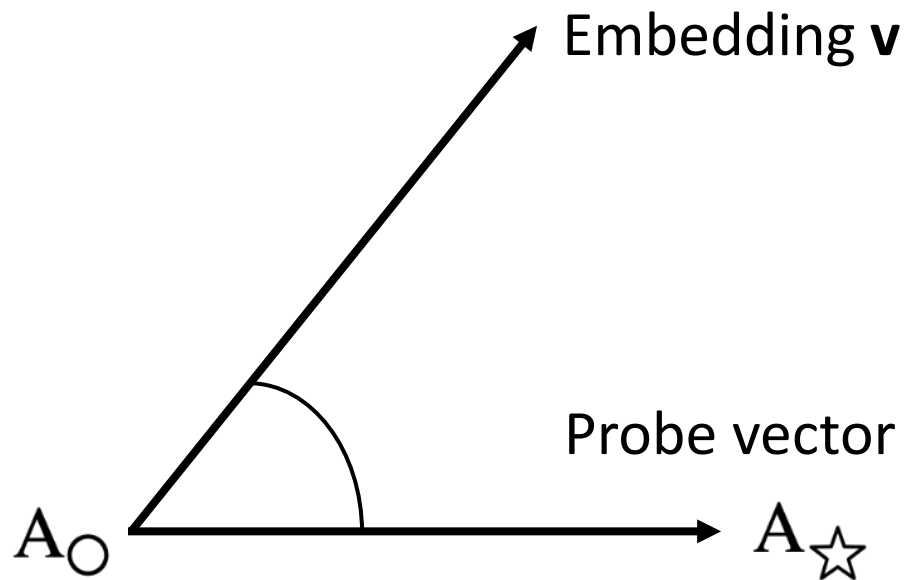
$A_{\star}$  = {*sharp, spiky, angular, jagged, hard, edgy, pointed, prickly, rugged, uneven*}

$A_{\circ}$  = {*round, circular, soft, fat, chubby, curved, smooth, plush, plump, rotund*}



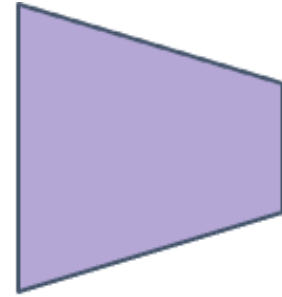
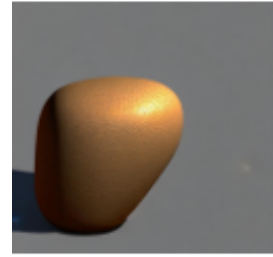
$$A_{\circ} \xrightarrow{\text{Probe vector}} A_{\star} \quad w_{adj} := \sum_{\langle a \rangle \in A_{\star}} \hat{w}_{\langle a \rangle} - \sum_{\langle a \rangle \in A_{\circ}} \hat{w}_{\langle a \rangle}$$

$A_{\star}$  = {*sharp, spiky, angular, jagged, hard, edgy, pointed, prickly, rugged, uneven*}  
 $A_{\circ}$  = {*round, circular, soft, fat, chubby, curved, smooth, plush, plump, rotund*}

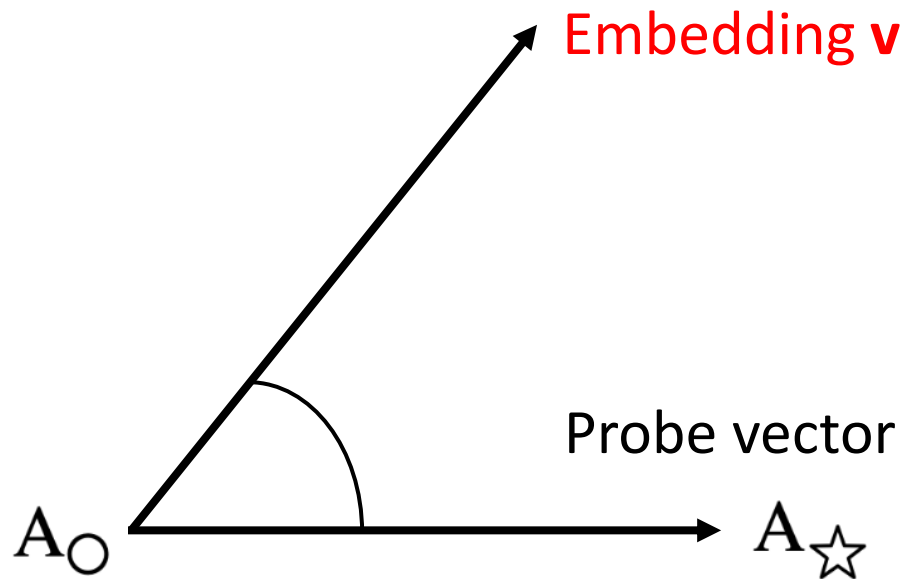


$$w_{adj} := \sum_{\langle a \rangle \in A_☆} \hat{w}_{\langle a \rangle} - \sum_{\langle a \rangle \in A_O} \hat{w}_{\langle a \rangle}$$

- $A_☆$  = {*sharp, spiky, angular, jagged, hard, edgy, pointed, prickly, rugged, uneven*}  
 $A_O$  = {*round, circular, soft, fat, chubby, curved, smooth, plush, plump, rotund*}



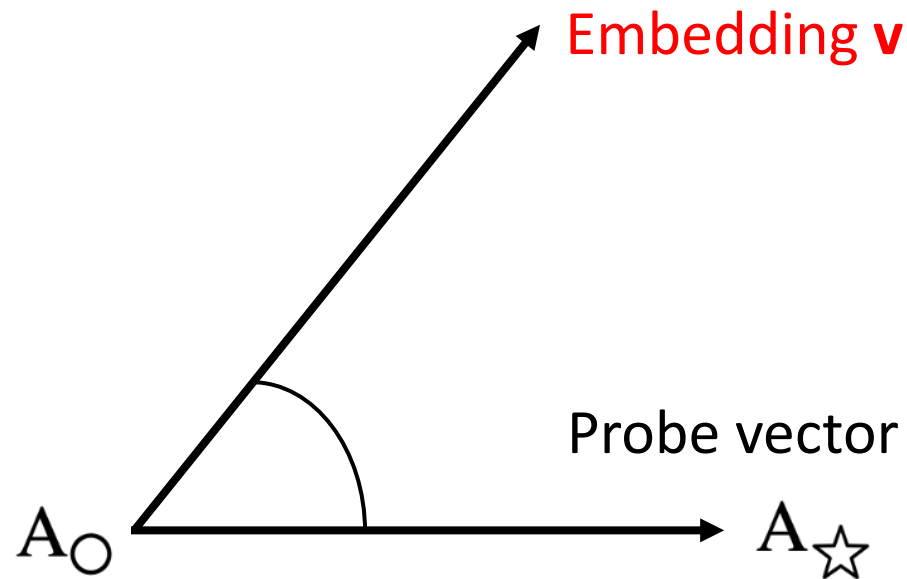
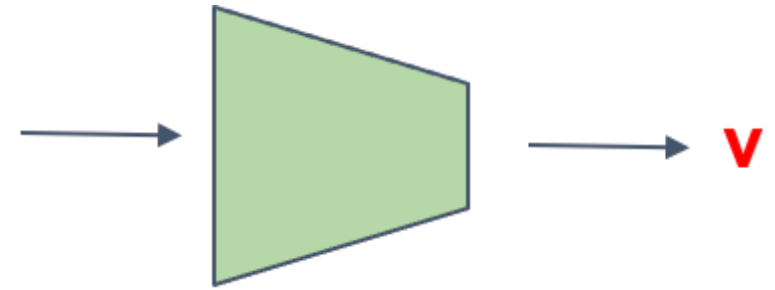
**v**



$$w_{adj} := \sum_{\langle a \rangle \in A_{\star}} \hat{w}_{\langle a \rangle} - \sum_{\langle a \rangle \in A_{\circ}} \hat{w}_{\langle a \rangle}$$

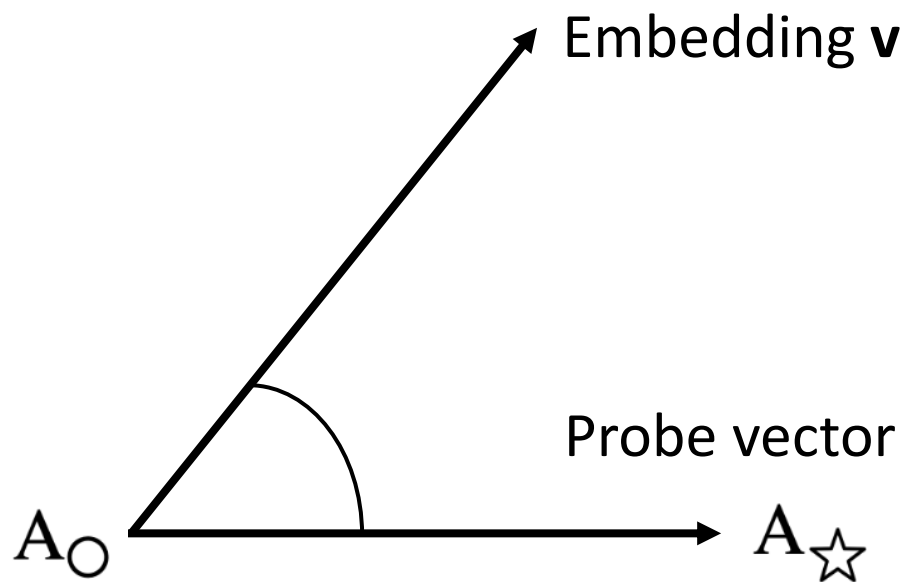
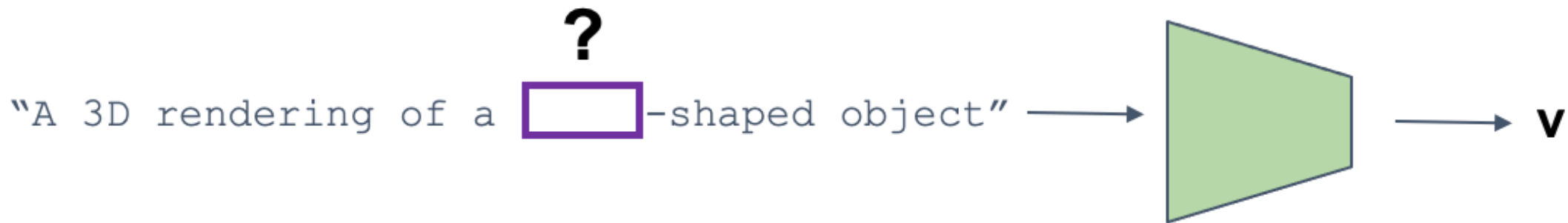
- $A_{\star}$  = {sharp, spiky, angular, jagged, hard, edgy, pointed, prickly, rugged, uneven}
- $A_{\circ}$  = {round, circular, soft, fat, chubby, curved, smooth, plush, plump, rotund}

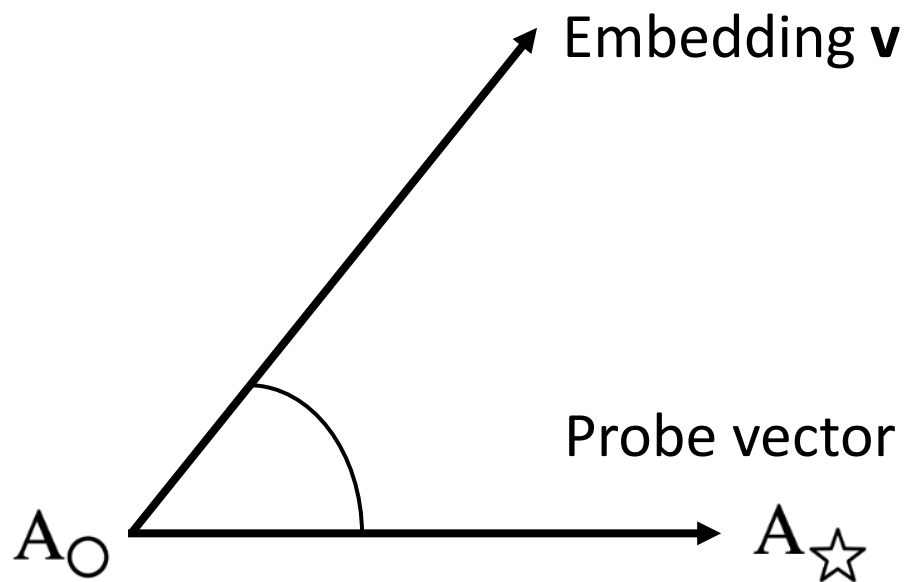
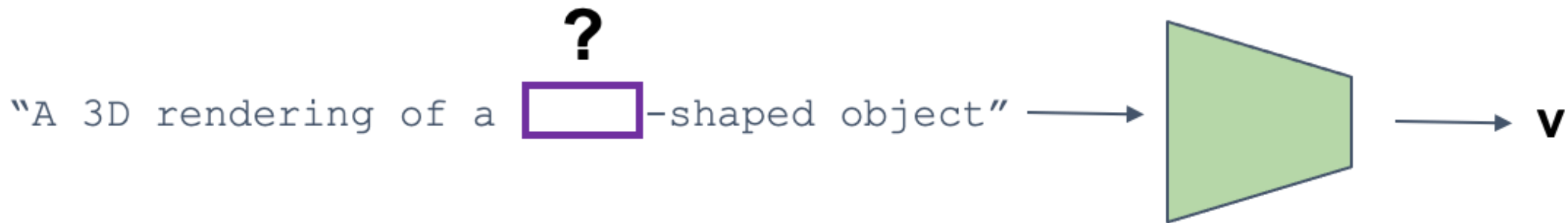
"A 3D rendering of a bouba-shaped object"



$$w_{adj} := \sum_{\langle a \rangle \in A_{\star}} \hat{w}_{\langle a \rangle} - \sum_{\langle a \rangle \in A_{\circ}} \hat{w}_{\langle a \rangle}$$

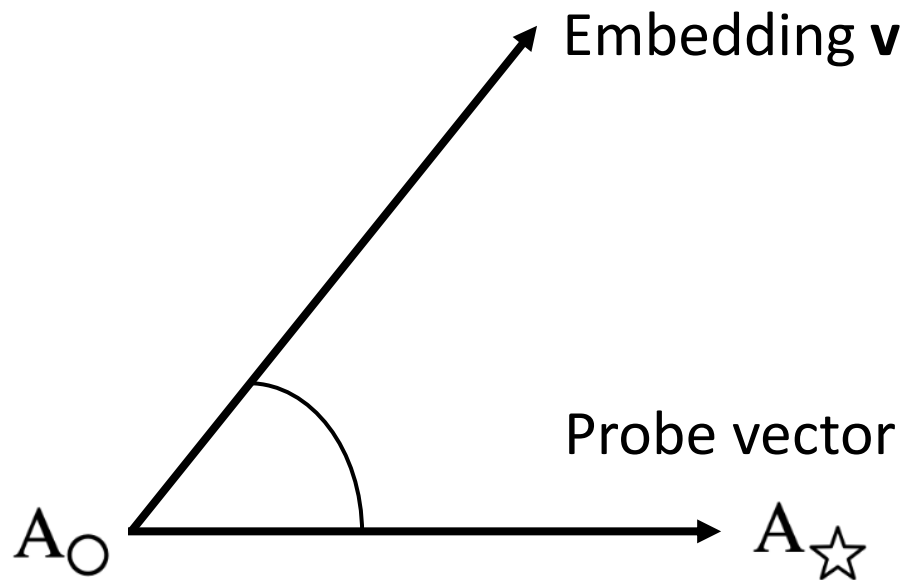
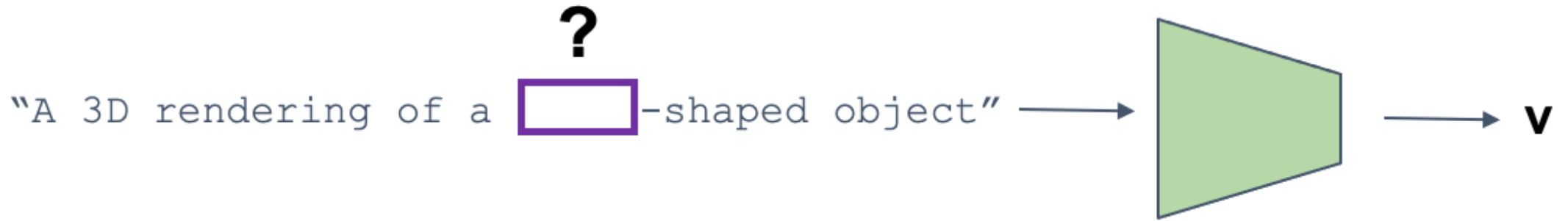
- $A_{\star}$  = {sharp, spiky, angular, jagged, hard, edgy, pointed, prickly, rugged, uneven}
- $A_{\circ}$  = {round, circular, soft, fat, chubby, curved, smooth, plush, plump, rotund}





$C_☆$ :	$\langle p \ t \ k \ s \ h \ x \rangle$
$V_☆$ :	$\langle e \ i \rangle$
$C_O$ :	$\langle b \ d \ g \ m \ n \ l \rangle$
$V_O$ :	$\langle o \ u \rangle$
$V_-$ :	$\langle a \rangle$

McCormick et al. (2015)



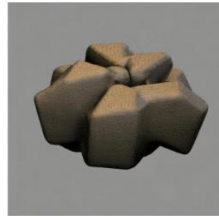
$C_{\star}$ :	$\langle p \ t \ k \ s \ h \ x \rangle$
$V_{\star}$ :	$\langle e \ i \rangle$
$C_{\circ}$ :	$\langle b \ d \ g \ m \ n \ l \rangle$
$V_{\circ}$ :	$\langle o \ u \rangle$
$V_{-}$ :	$\langle a \rangle$

McCormick et al. (2015)

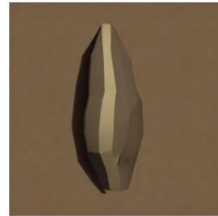
$\Psi_{\star}$ :	<i>kitaki</i>	<i>hatiha</i>	<i>pepape</i>	<i>xisixi</i>	<i>hipehi</i>	<i>xaxaxa</i>	<i>texete</i>	...
$\Psi_{\circ}$ :	<i>gugagu</i>	<i>bodubo</i>	<i>gunogu</i>	<i>daluda</i>	<i>momomo</i>	<i>lunulu</i>	<i>gadaga</i>	...

# Pseudowords sorted by geometric score

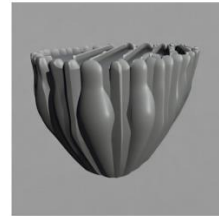
High ( ☆ )



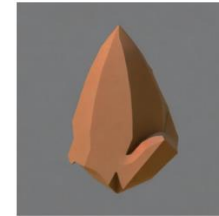
*tepite* ☆



*kataka* ☆



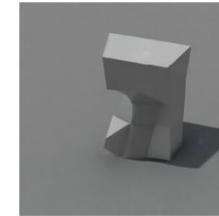
*gugugu* ○



*pitapi* ☆

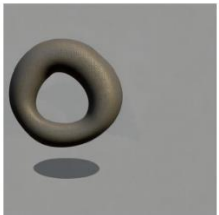


*daguda* ○

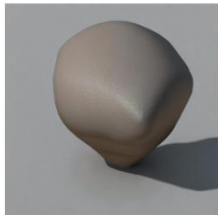


*setese* ☆

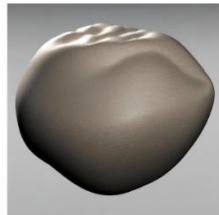
Low ( ○ )



*nubanu* ○



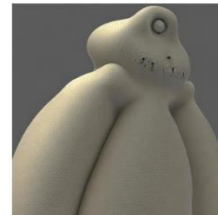
*gomago* ○



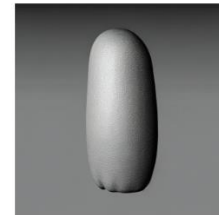
*magama* ○



*bomabo* ○



*mamoma* ○

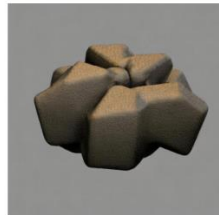


*lobalo* ○

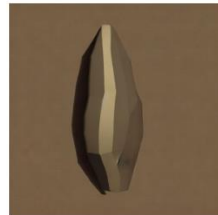


## Pseudowords sorted by geometric score

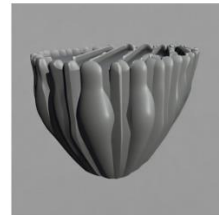
High ( ☆ )



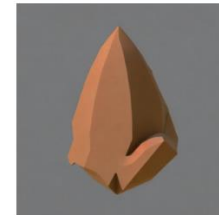
*tepite* ☆



*kataka* ☆

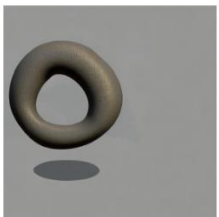


*gugugu* ○

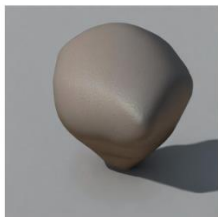


*pitapi* ☆

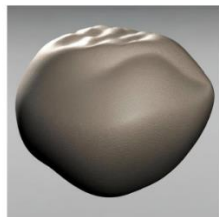
Low ( ○ )



*nubanu* ○



*gomago* ○



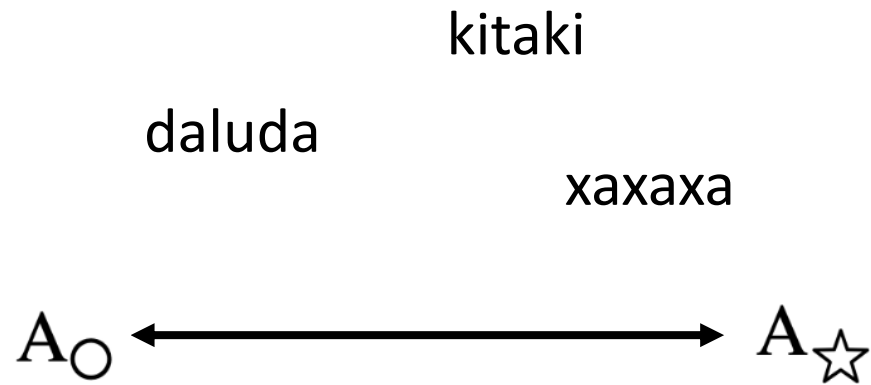
*magama* ○



*bomabo* ○

Model	AUC	$\tau$
Stable Diffusion	0.74	0.34
CLIP	0.77	0.39
(random)	0.50	0.00

# Geometric scoring



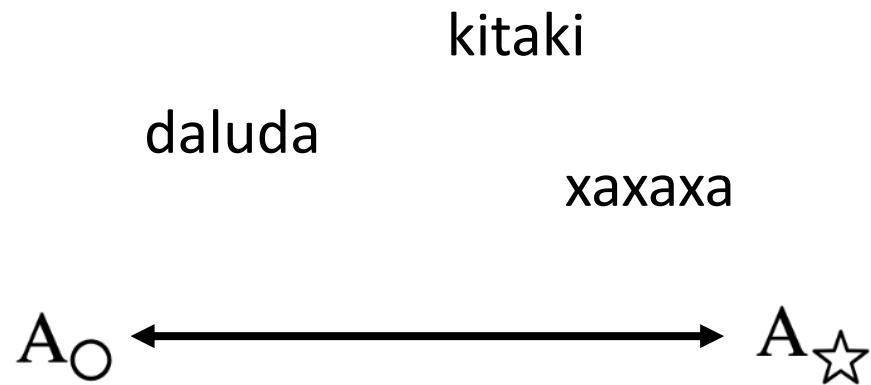
$\Psi_{\star}$ : *kitaki*    *hatiha*    *pepape*    *xisixi*    *hipehi*    *xaxaxa*    *texete*    ...

$\Psi_{\circ}$ : *gugagu*    *bodubo*    *gunogu*    *daluda*    *momomo*    *lunulu*    *gadaga*    ...

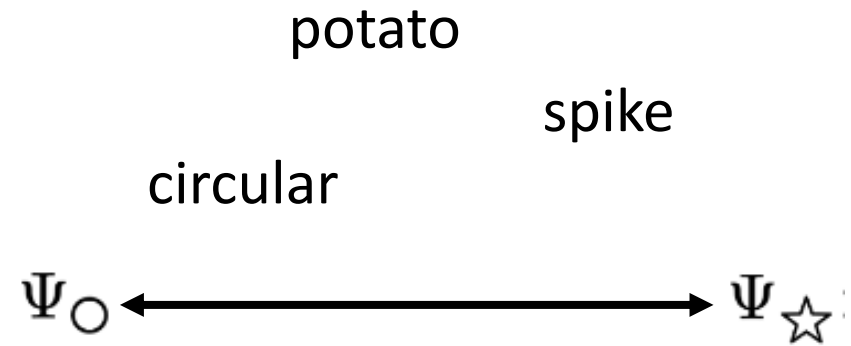
$A_{\star}$  = {*sharp, spiky, angular, jagged, hard, edgy, pointed, prickly, rugged, uneven*}

$A_{\circ}$  = {*round, circular, soft, fat, chubby, curved, smooth, plush, plump, rotund*}

## Geometric scoring



## Phonetic scoring



$\Psi_{\star}$ :    *kitaki    hatiha    pepape    xisixi    hipehi    xaxaxa    texete    ...*

$\Psi_{\circ}$ :    *gugagu    bodubo    gunogu    daluda    momomo    lunulu    gadaga    ...*

$A_{\star}$     =    *{sharp, spiky, angular, jagged, hard, edgy, pointed, prickly, rugged, uneven}*

$A_{\circ}$     =    *{round, circular, soft, fat, chubby, curved, smooth, plush, plump, rotund}*

# Real English words sorted by phonetic score

## Low ( ○ )

## High ( ☆ )

### Stable Diffusion

Noun *butterball, yolk, pregnancy, booger, eggnog, turnip, bellyful, crybaby, doughboy*

*shard, kite, origami, hexagon, diamond, flake, octagon, triangle, protractor, lozenge, foldout*

Adj. *obese, chubby, stinky, pudgy, overweight, fat, pregnant, plump, drowsy, soggy, squishy*

*triangular, diagonal, angular, shattering, jagged, rectangular, edgy, housebroken, geometrical*

### CLIP

Noun *doughboy, loudmouth, gumdrop, boogeyman, madwoman, lord, butterball, goddaughter*

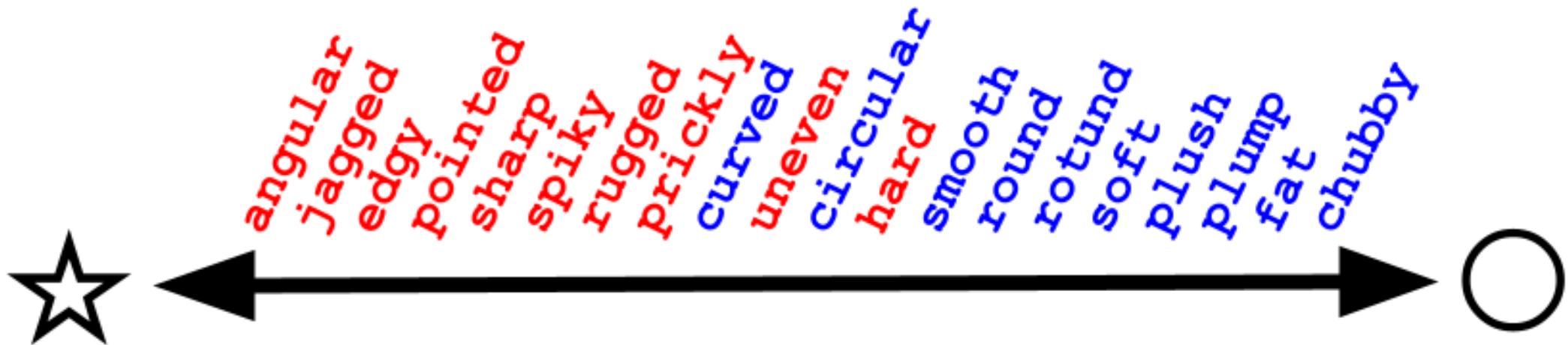
*prefix, talkativeness, asker, flexibility, shears, shift, peek, slope, task, exit, hemline, tightness*

Adj. *muggy, soggy, gloomy, grouchy, lumpy humongous, hungry, cloudless, unsmiling*

*apelike, flexible, diagonal, static, triangular, external, shipshape, interlocking, angular*

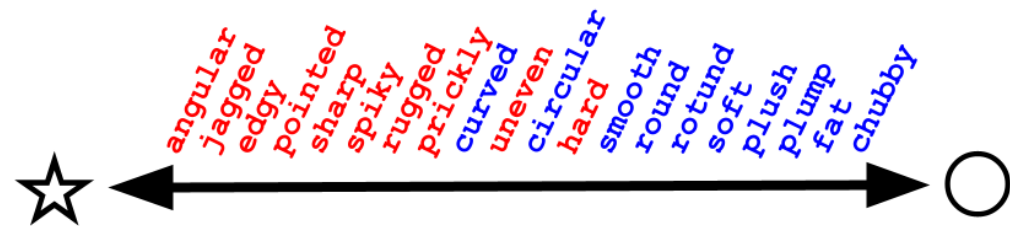
# GT adjectives sorted by phonetic score

## Stable Diffusion

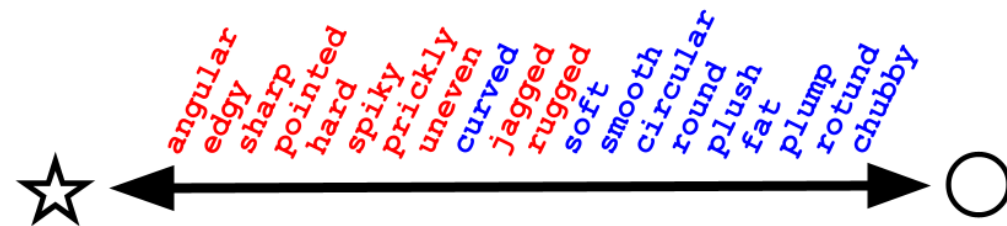


# GT adjectives sorted by phonetic score

Stable Diffusion

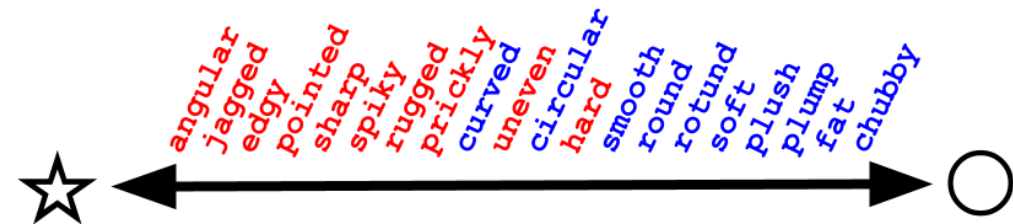


CLIP

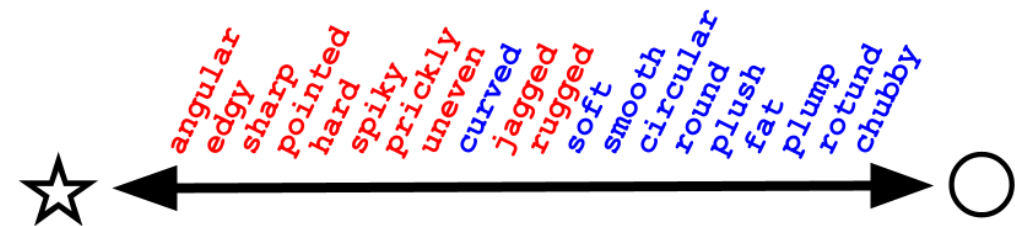


# GT adjectives sorted by phonetic score

## Stable Diffusion



## CLIP

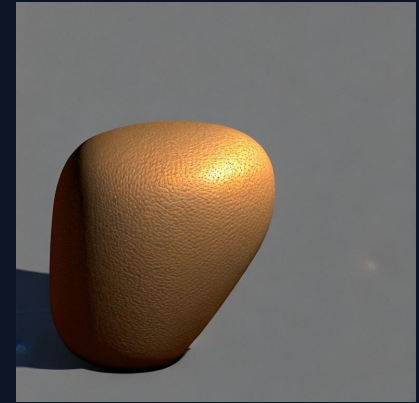


Model	AUC	$\tau$
Stable Diffusion	0.97	0.68
CLIP	0.98	0.70
(random)	0.50	0.00

## Conclusion

We show that VLMs have learned sound symbolism, providing:

- **Insight and interpretability** for black-box models
- **A new computational tool** for studying sound symbolism in language





# Questions?

