

# Reflection-based Word Attribute Transfer

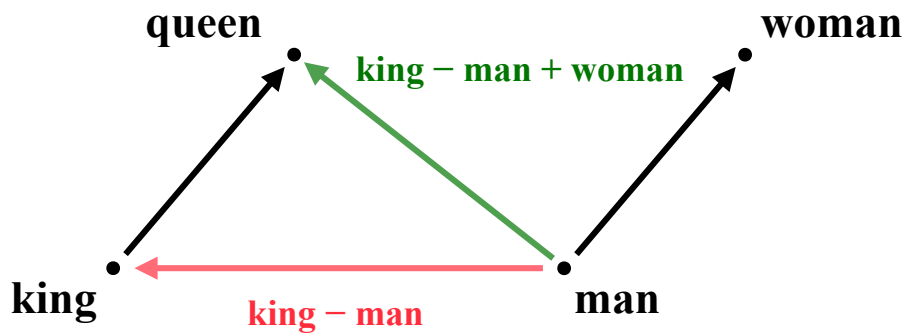
Yoichi Ishibashi, Katsuhito Sudoh, Koichiro Yoshino, Satoshi Nakamura

Nara Institute of Science and Technology (NAIST), Japan

## Analogy in the embedding space

- is a operation that transfer word attributes
- Change word attributes (e.g. gender)

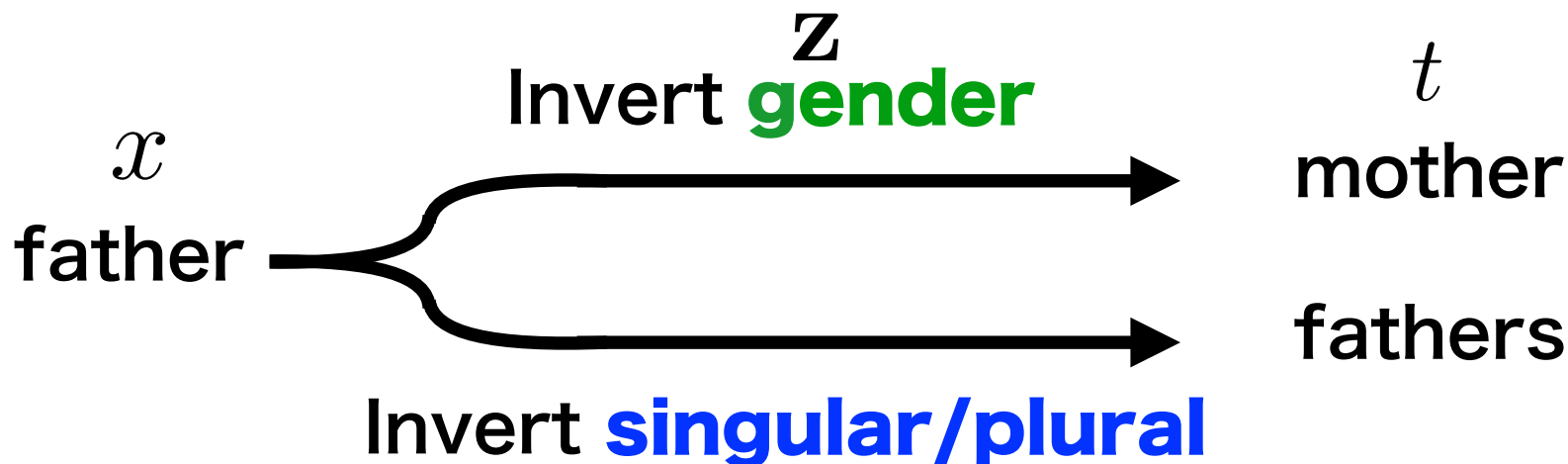
$$\overrightarrow{king} - \overrightarrow{man} + \overrightarrow{woman} \approx \overrightarrow{queen}$$



## Word attribute transfer task

- Get a word vector that inverted attribute of an input word vector

$$\mathbf{v}_{\text{mother}} \approx f_{\text{gender}}(\mathbf{v}_{\text{father}})$$



What can word attribute transfer be used for?

- E.g. Data Augmentation

Target attribute	Input words	Output words
Gender	I am <b>his</b> <b>mother</b> .	I am <b>her</b> <b>father</b> .
Antonym	<b>Nobody</b> has a suit.	<b>Someone</b> has a suit.
Capital-Country	I live in <b>Japan</b> .	I live in <b>Tokyo</b> .

$f_z$  transfer words if they have a target attribute  $Z$

- E.g. **man**  $\rightarrow$  **woman** (attribute: gender)

$$\mathbf{v}_{\text{woman}} \approx f_{\text{gender}}(\mathbf{v}_{\text{man}})$$

$f_z$  does not transfer words if it does not has a target attribute  $Z$

- E.g. **person**  $\rightarrow$  **person** (attribute: gender)

$$\mathbf{v}_{\text{person}} \approx f_{\text{gender}}(\mathbf{v}_{\text{person}})$$

# Analogy-based Word Attribute Transfer 6

## Analogy-based word attribute transfer

king  $-$  (man  $-$  woman) = queen

queen  $+$  (man  $-$  woman) = king

- **Add** or **subtract** a difference vector

## Problem

- **Need explicit knowledge** whether input word has the target attribute or not

## Goal

- Transform word attributes **without the explicit knowledge**

**Proposed method**

## What is an **ideal transfer function**?

- No explicit knowledge  
= Transfer any words with the same function

$$\mathbf{v}_{\text{man}} = f(\mathbf{v}_{\text{woman}})$$

$$\mathbf{v}_{\text{woman}} = f(\mathbf{v}_{\text{man}})$$

$$\mathbf{v}_{\text{person}} = f(\mathbf{v}_{\text{person}})$$

Combine above formulas

$$\mathbf{v}_x = f(f(\mathbf{v}_x))$$

← **Nature of the  
ideal function**



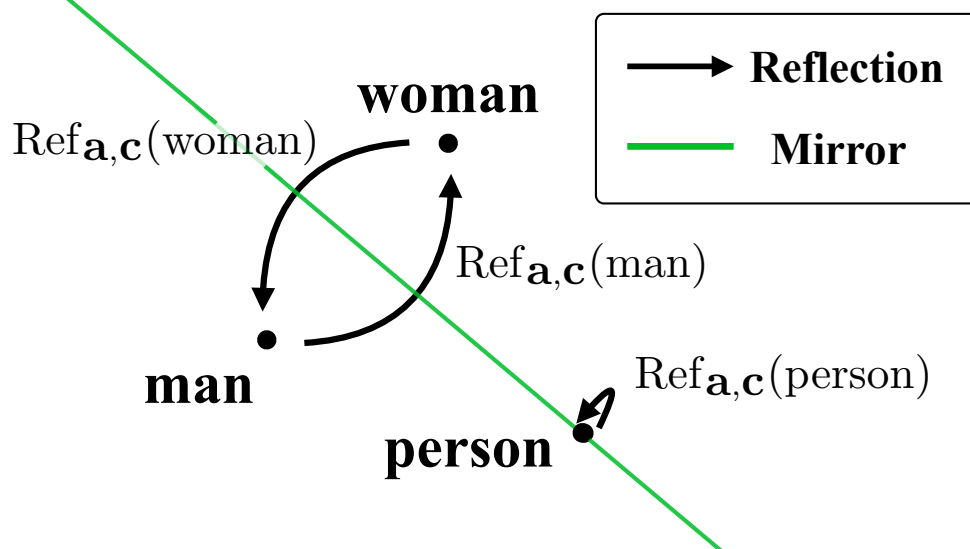
## Reflection is an ideal function

- Transfer any words with the same function

$$\mathbf{v} = \text{Ref}_{\mathbf{a},\mathbf{c}}(\text{Ref}_{\mathbf{a},\mathbf{c}}(\mathbf{v})) \quad \text{Ref}_{\mathbf{a},\mathbf{c}}(\mathbf{v}) = \mathbf{v} - 2 \frac{(\mathbf{v} - \mathbf{c}) \cdot \mathbf{a}}{\mathbf{a} \cdot \mathbf{a}} \mathbf{a}$$

- Move two vectors through a hyperplane (**mirror**)

man =  $\text{Ref}_{\mathbf{a},\mathbf{c}}(\text{woman})$   
woman =  $\text{Ref}_{\mathbf{a},\mathbf{c}}(\text{man})$   
person =  $\text{Ref}_{\mathbf{a},\mathbf{c}}(\text{person})$



# Reflection-based Word Attribute Transfer 10

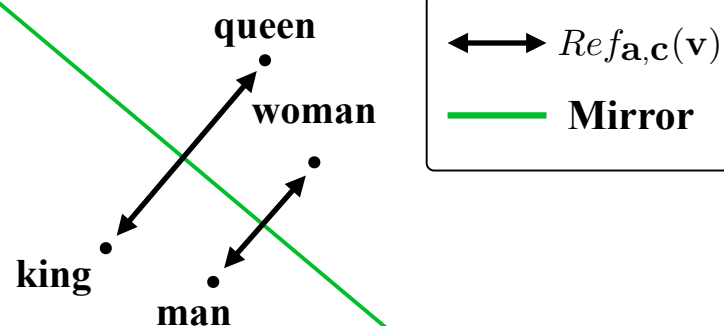
## How to apply to word attribute transfer?

- Transfer an input word vector  $\mathbf{v}_x$  to a target word vector  $\mathbf{v}_t$

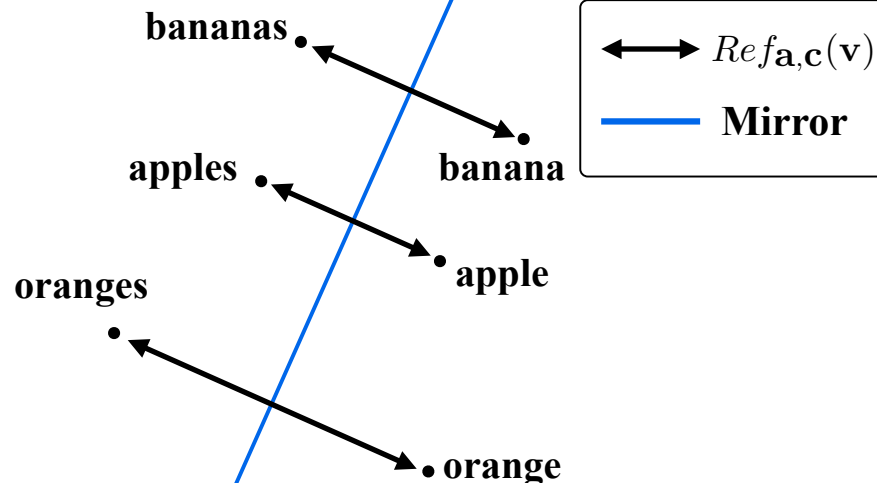
$$\mathbf{v}_t \approx \mathbf{v}_y = \text{Ref}_{\mathbf{a}, \mathbf{c}}(\mathbf{v}_x)$$

- Learn a mirror for each attributes

Male $\leftrightarrow$ Female



Singular $\leftrightarrow$ Plural



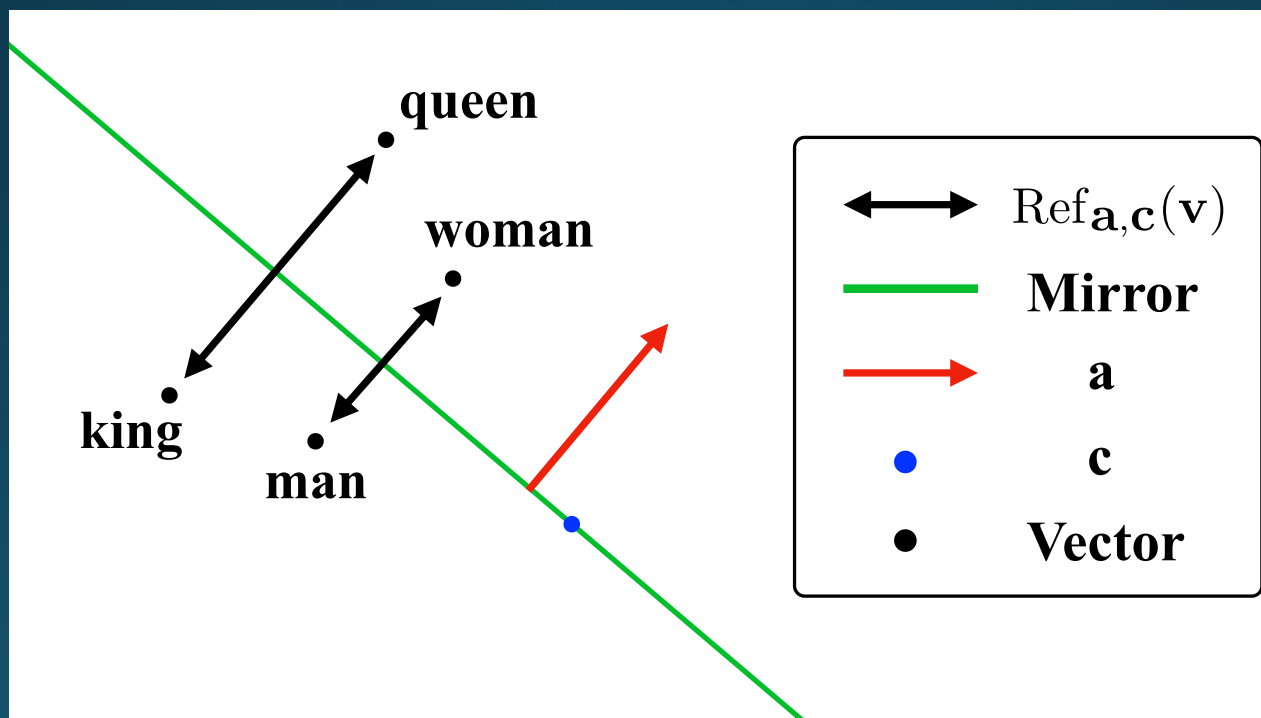
# Reflection-based Word Attribute Transfer 11

How to learn the mirror?

Idea : **Estimate** **a** **and** **c** **by** **MLP**

**a** ... A vector orthogonal to the mirror

**c** ... A point through which the mirror passes



# Two types of mirror estimation

12

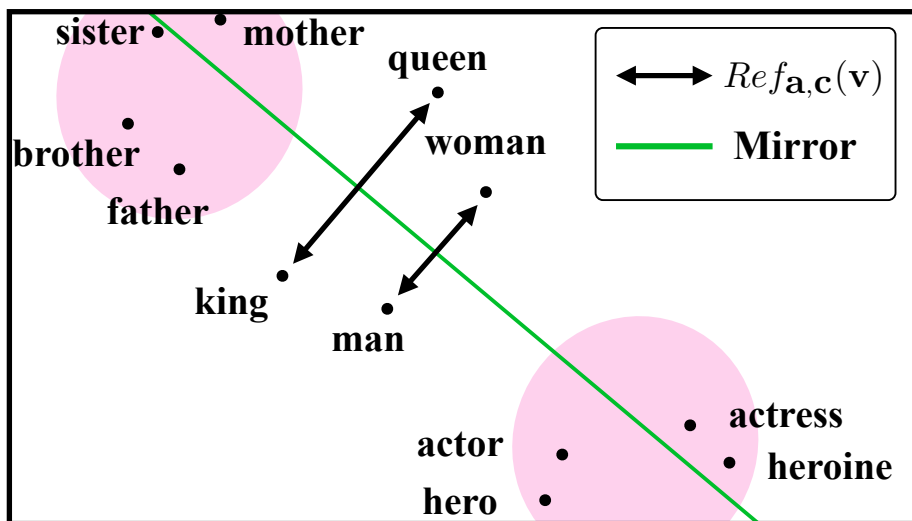
## ① Single mirror

Estimate from an attribute  $\mathbf{z}$

⇒ Some pairs are  
**non-transferable**

$$\mathbf{a} = \text{MLP}_{\theta_1}(\mathbf{z})$$

$$\mathbf{c} = \text{MLP}_{\theta_2}(\mathbf{z})$$



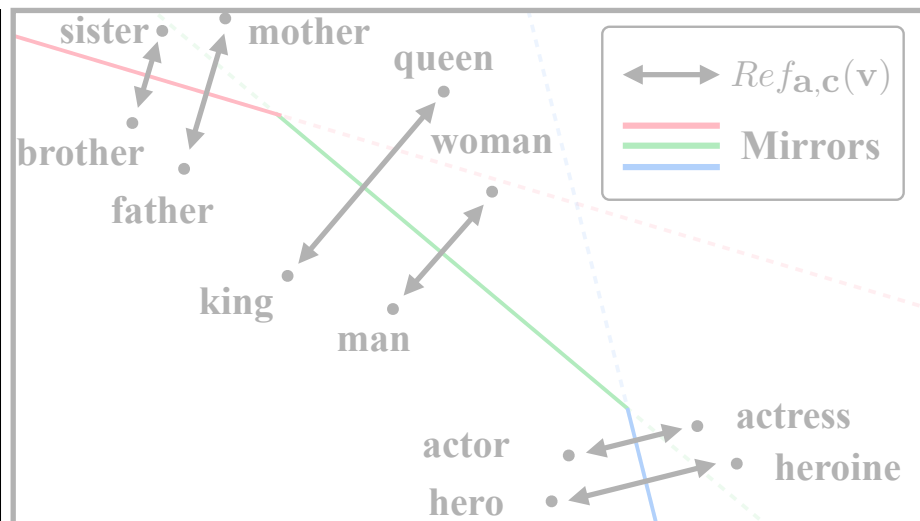
## ② Parameterized mirrors

Estimate from  $\mathbf{z}$

and an input word vector  $\mathbf{v}_x$   
⇒ Work more flexibly

$$\mathbf{a} = \text{MLP}_{\theta_1}([\mathbf{z}; \mathbf{v}_x])$$

$$\mathbf{c} = \text{MLP}_{\theta_2}([\mathbf{z}; \mathbf{v}_x])$$



# Two types of mirror estimation

13

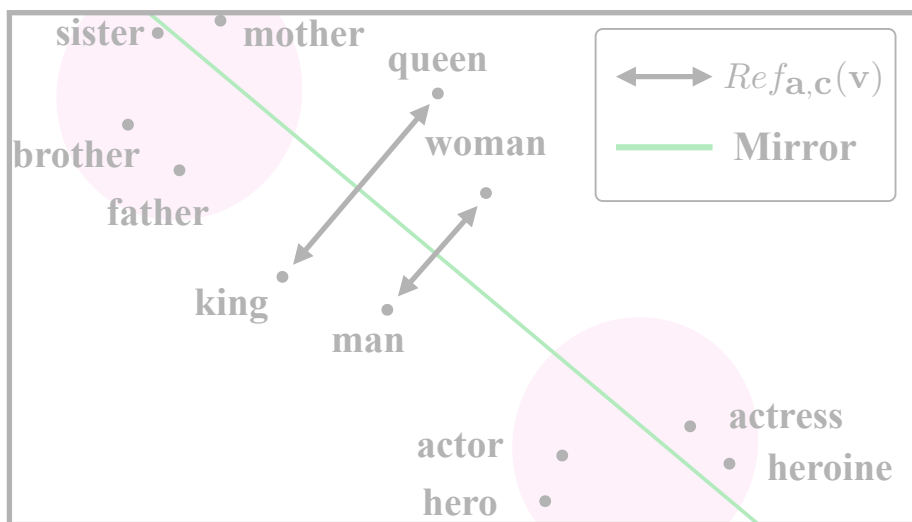
## ① Single mirror

Estimate from an attribute  $\mathbf{z}$

⇒ Some pairs are non-transferable

$$\mathbf{a} = \text{MLP}_{\theta_1}(\mathbf{z})$$

$$\mathbf{c} = \text{MLP}_{\theta_2}(\mathbf{z})$$



## ② Parameterized mirrors

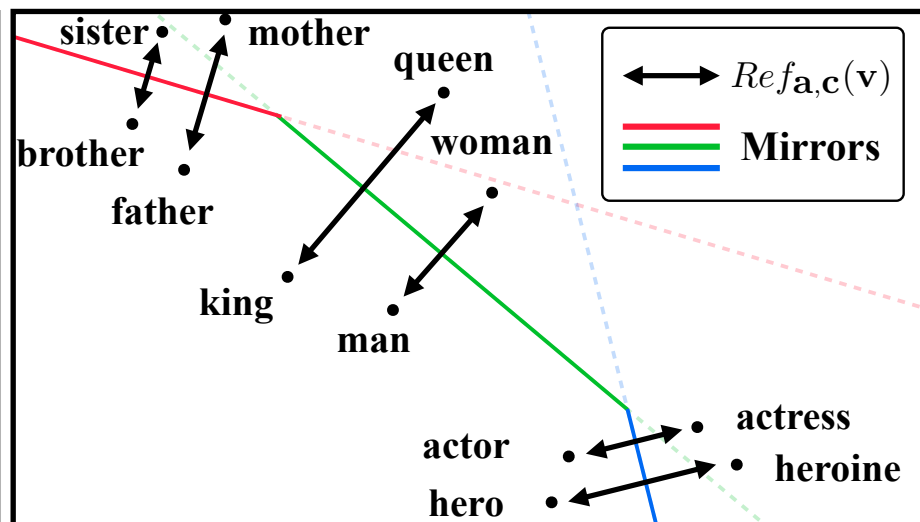
Estimate from  $\mathbf{z}$

and an input word vector  $\mathbf{v}_x$

⇒ **Work more flexibly**

$$\mathbf{a} = \text{MLP}_{\theta_1}([\mathbf{z}; \mathbf{v}_x])$$

$$\mathbf{c} = \text{MLP}_{\theta_2}([\mathbf{z}; \mathbf{v}_x])$$



# Experiments

Compare reflection and baselines

- **Four different attributes**
  - Male-Female, Singular-Plural  
Capital-Country, Antonyms
- **Two pre-trained word embeddings**
  - word2vec (SGNS), GloVe
- **Two evaluation metrics**
  - Accuracy, Stability

- Attribute words ... Four different binary attributes

Attribute (z)	Train	Val	Test	Example (x, t)
<b>Male-Female (MF)</b>	29	12	12	(king, queen)
<b>Singular-Plural (SP)</b>	90	25	25	(king, kings)
<b>Capital-Country (CC)</b>	59	25	25	(Japan, Tokyo)
<b>Antonym (AN)</b>	1354	290	290	(good, bad)

- Non-attribute words

- Train  $0 \leq |\mathcal{N}_{train}| \leq 50$

Test  $|\mathcal{N}_{test}| = 1000$



① **Accuracy:** Ratio of attribute words transferred



② **Stability:** Ratio of **non**-attribute words **not** transferred



**Best method:** Reflection with parameterized mirrors  
→ High performance in both accuracy and stability

**Worst method:** MLP

Method	GloVe							
	Accuracy (%)				Stability (%)			
	MF	SP	CC	AN	MF	SP	CC	AN
Ref	12.5	2.0	26.0	0.0	100.0	100.0	100.0	100.0
Ref+PM	<b>45.8</b>	<b>50.0</b>	<b>76.0</b>	33.5	99.7	99.1	99.2	100.0
MLP	4.2	10.0	18.0	<b>36.7</b>	5.1	7.0	5.2	1.2
Diff+	25.0	2.0	26.0	-	99.3	94.2	99.3	-
Diff-	25.0	2.0	24.0	-	100.0	99.9	99.5	-

MF: Male-Female, SP: Singular-Plural, CC: Country-Capital, AN: Antonym

**Best method:** Reflection with parameterized mirrors  
→ High performance in both accuracy and stability

**Worst method:** MLP

Method	GloVe							
	Accuracy (%)				Stability (%)			
	MF	SP	CC	AN	MF	SP	CC	AN
Ref	12.5	2.0	26.0	0.0	100.0	100.0	100.0	100.0
Ref+PM	45.8	50.0	76.0	33.5	99.7	99.1	99.2	100.0
MLP	4.2	10.0	18.0	36.7	5.1	7.0	5.2	1.2
Diff+	25.0	2.0	26.0	-	99.3	94.2	99.3	-
Diff-	25.0	2.0	24.0	-	100.0	99.9	99.5	-

MF: Male-Female, SP: Singular-Plural, CC: Country-Capital, AN: Antonym

Reflection with parameterized mirrors (Ref+PM) can **selectively** transfer words without the knowledge

boy  girl

you  you

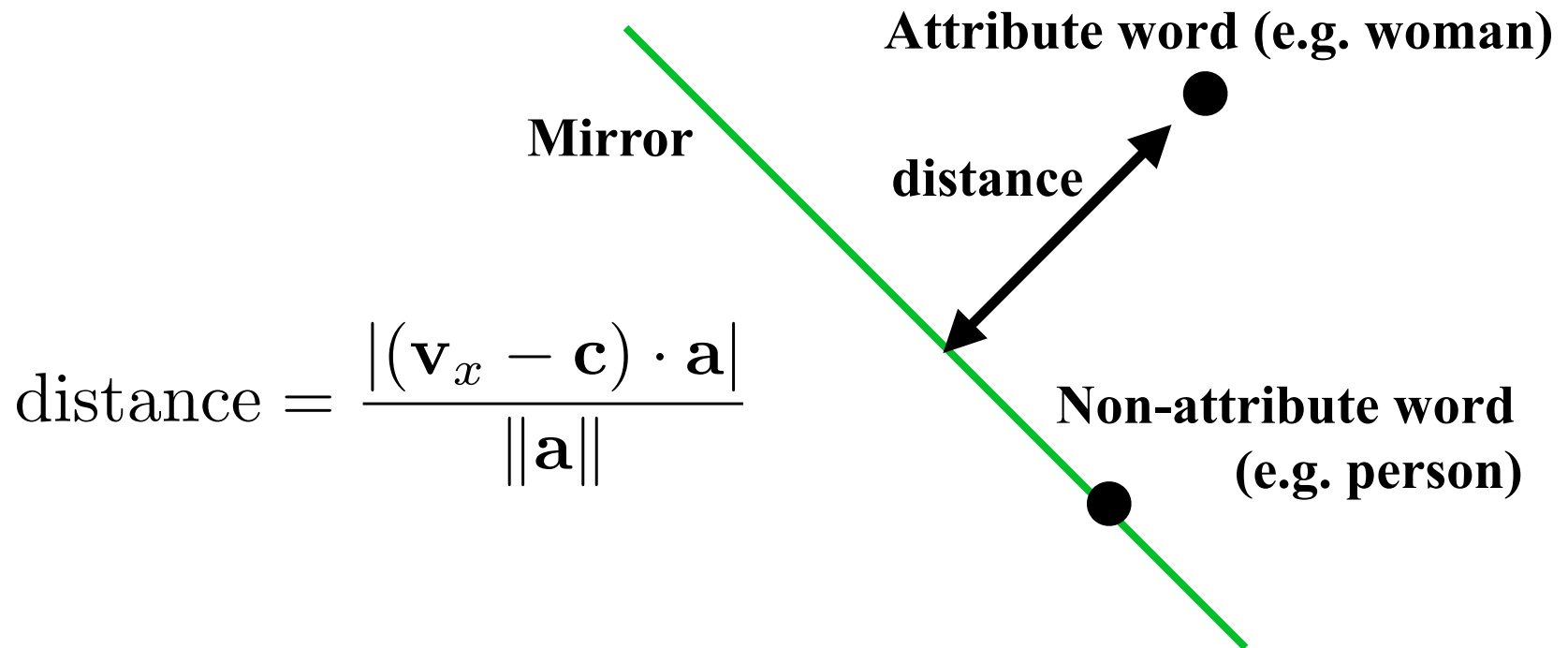
Input	the <b>woman</b> got married when you were a <b>boy</b> .
Ref	the <b>woman</b> got married when you were a <b>boy</b> .
Ref+PM	the <b>man</b> got married when you were a <b>girl</b> .
MLP	By_Katie_Klingsporn <b>girlfriend</b> Valerie_Glodowski fiancee Doughty_Evening_Chronicle ma'am Bob_Grossweiner_& a <b>mother</b> .
Diff+	the <b>man</b> got married when you were a <b>boy</b> .
Diff-	the <b>woman</b> got married when you were a <b>girl</b> .

# Why is the reflection very stable?

21

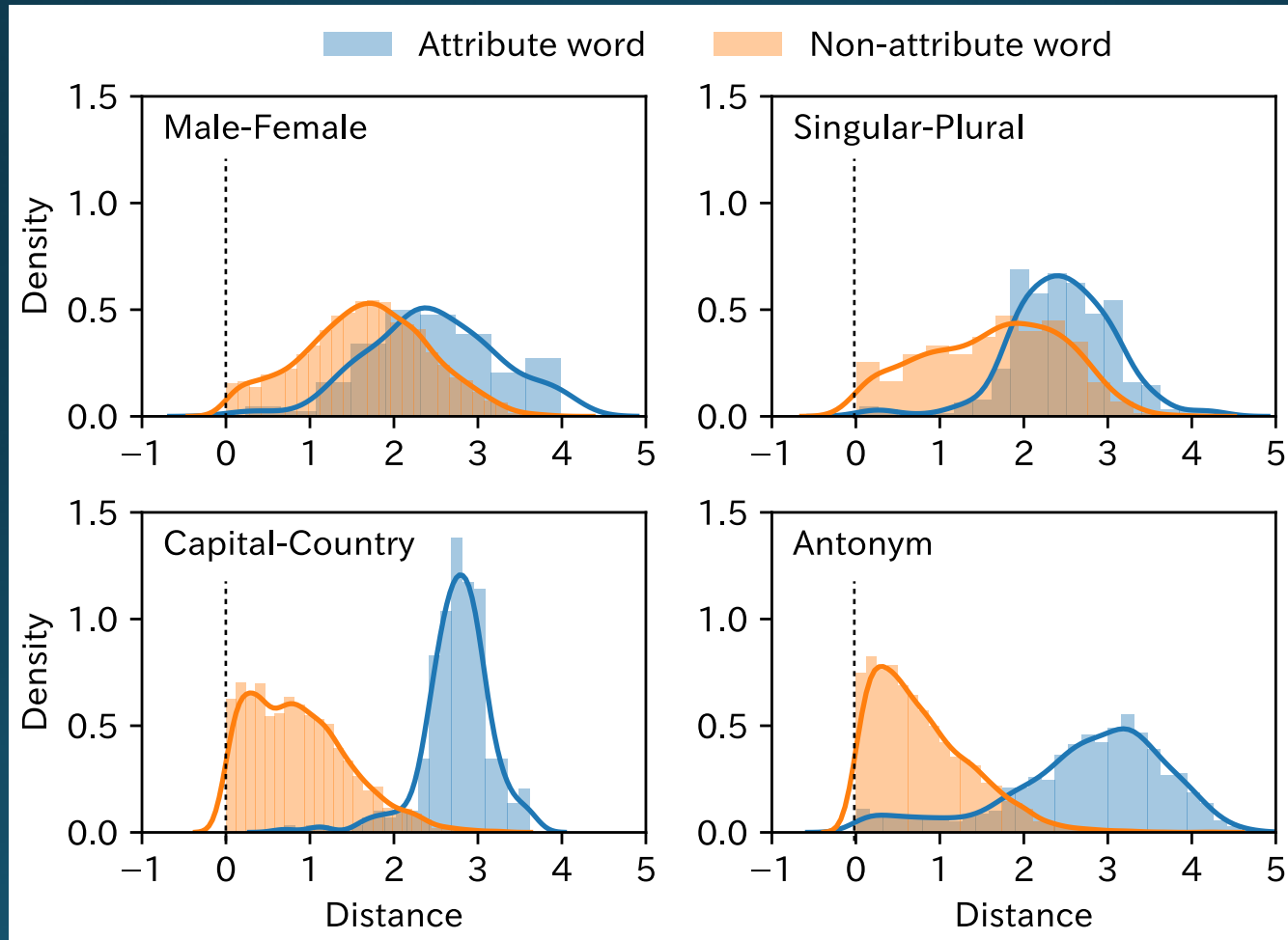
Hypothesis: **Non-attribute word distributes on its mirror**

→ Visualize the distance between a word vector and its mirror



# Distance between the word and its mirror 22

- Attribute words distributed **apart from the mirror**
- Non-attribute words distributed **near the mirror**



# Summary

## Background

- Word attribute transfer task
- Analogy can be used for the transfer
- Analogy-based transfer requires the explicit knowledge

## Proposed method

- Reflection-based word attribute transfer
- Reflection is an ideal mapping for word attribute transfer

## Experimental results

- Reflection-based transfer achieved best performance
- Reflection transfers attribute words  
e.g. man  $\rightarrow$  woman
- Reflection does not transfer non-attribute words  
e.g. person  $\rightarrow$  person