

Text Style Transfer with Confounders

Tianxiao Shen Regina Barzilay Tommi Jaakkola

tianxiao@mit.edu

What Is “Style Transfer”?

source

target



Monet → photo



horse → zebra

[Zhu et al. 2017]

From informal to formal

Gotta see both sides of the story

→ You have to consider both sides of the story

[Rao et al. 2018]

From Shakespeare to modern

Send thy man away → Send your man away

[Xu et al. 2012]

From negative to positive sentiment

I would recommend find another place.

→ I would recommend this place again!

[Shen et al. 2017]

From dialect to written standard

From complex to simple sentences

...

Easy: Paired Training Sets

- Supervised learning using paired examples of style transfer

source
(e.g., negative reviews)

owner: a very rude man.

i would not recommend giving them a try!

we were both so disappointed!

consistently slow.



target
(e.g., positive reviews)

owner: a very friendly man.

i'd definitely recommend giving them a try!

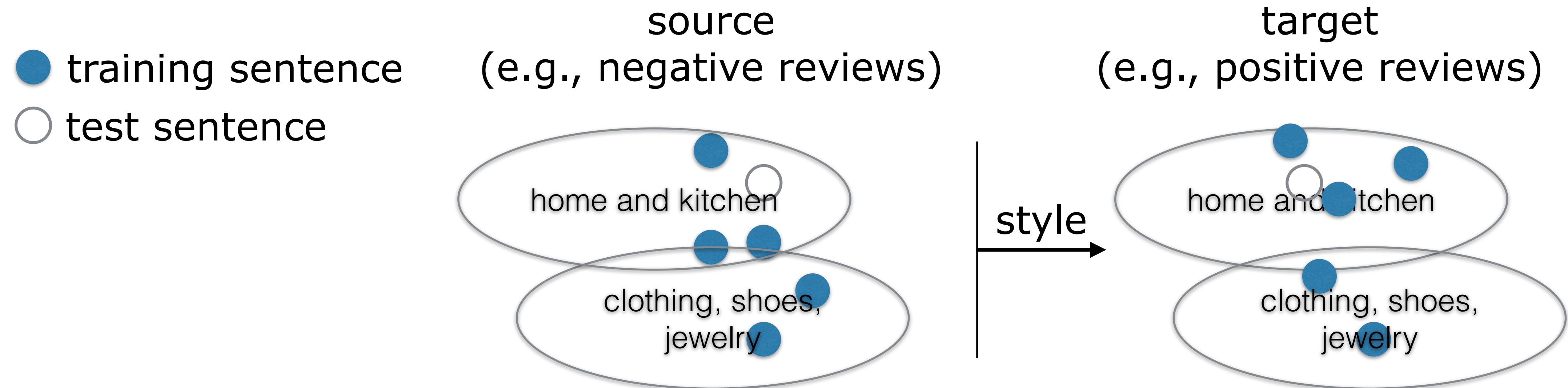
we were both so impressed!

consistently fast.

- To collect parallel data is very costly or even impossible

Intermediate: Unpaired but Distributionally Matched Sets

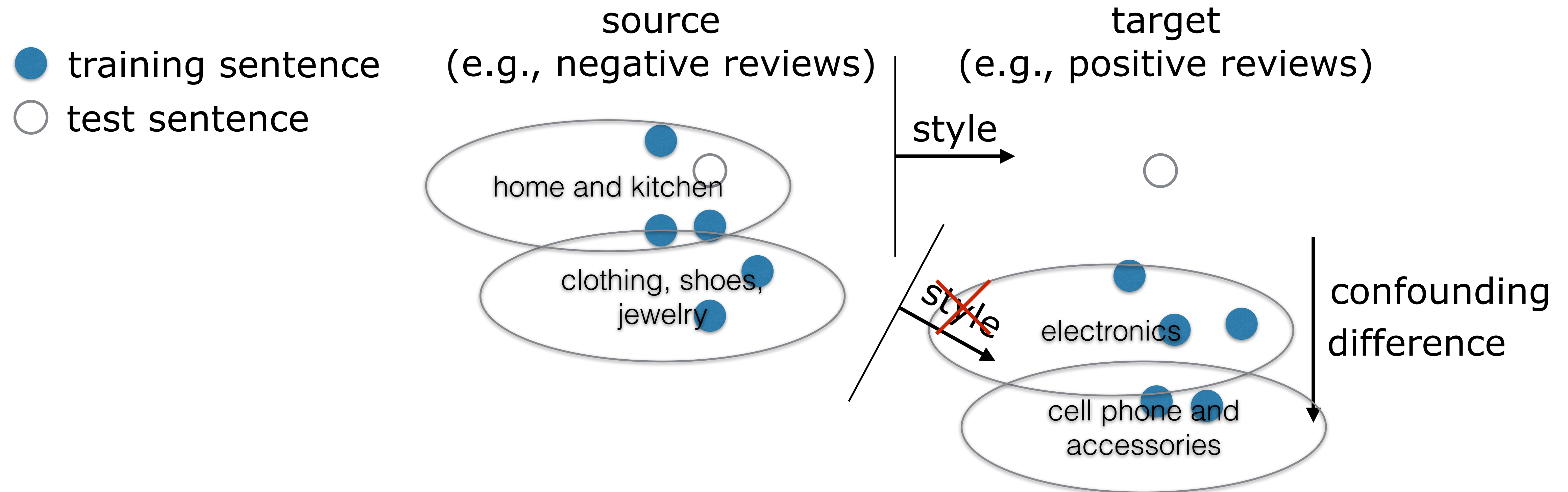
- Available source and target sentences as sets differ only in terms of style, i.e., they are distributionally matched otherwise



- The desired style change is just the source vs target difference
- New sentences map to sentences similar to those already seen during training

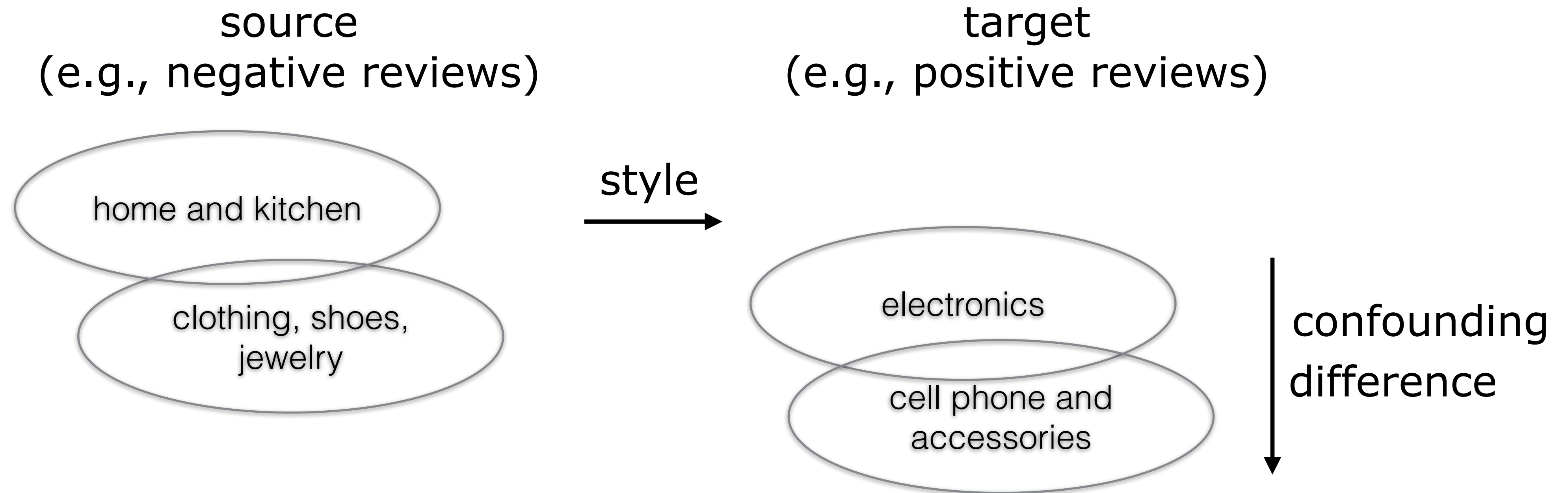
Hard: Unpaired, Not Distributionally Matched Sets

- There are **additional confounding differences** between source and target sentences



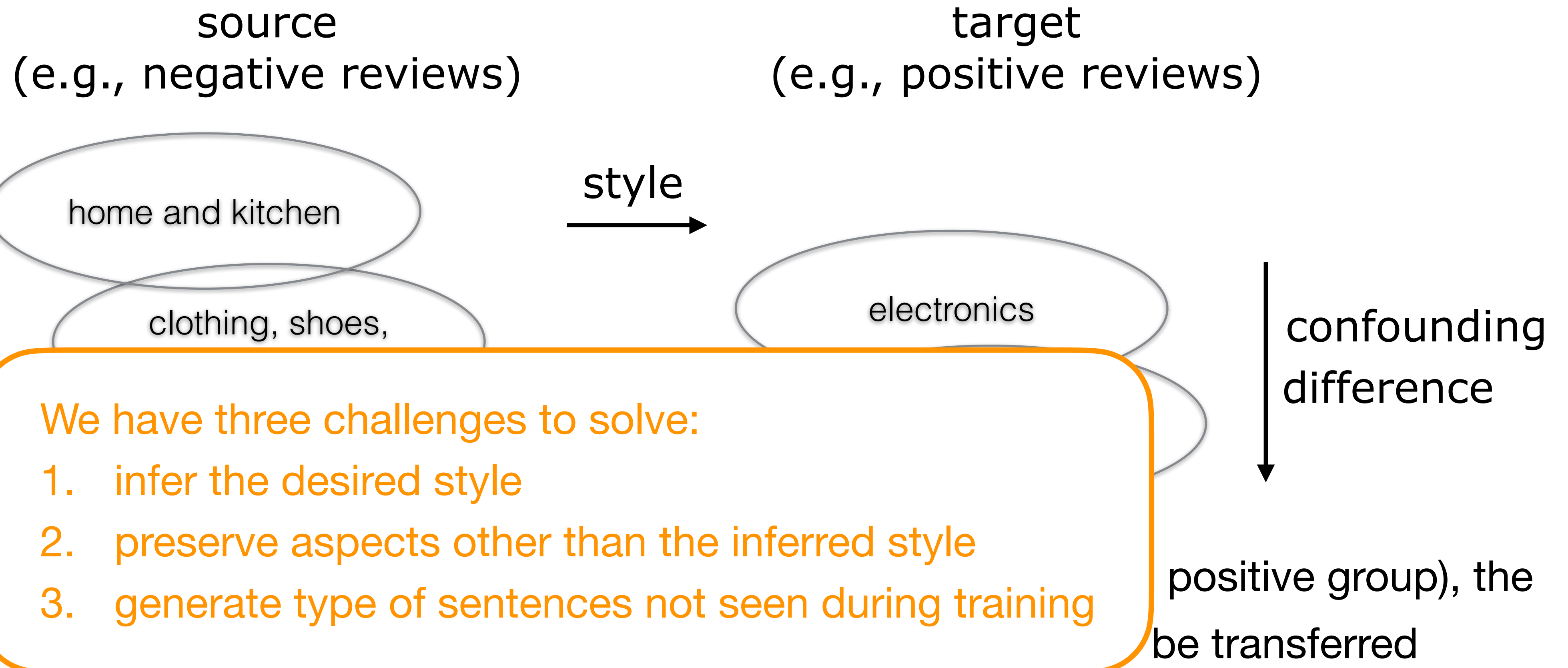
- Style change no longer equals source vs target difference
- New sentences map to type of sentences not seen during training

Solving Style Transfer with Confounders



- The task is illustrated by two groups of datasets (negative group and positive group), the **primary distinction** between them (sentiment) specifies the style to be transferred
- The intra-group variations (category) are **confounding differences** which need to be differentiated from the style and preserved during transfer

Solving Style Transfer with Confounders



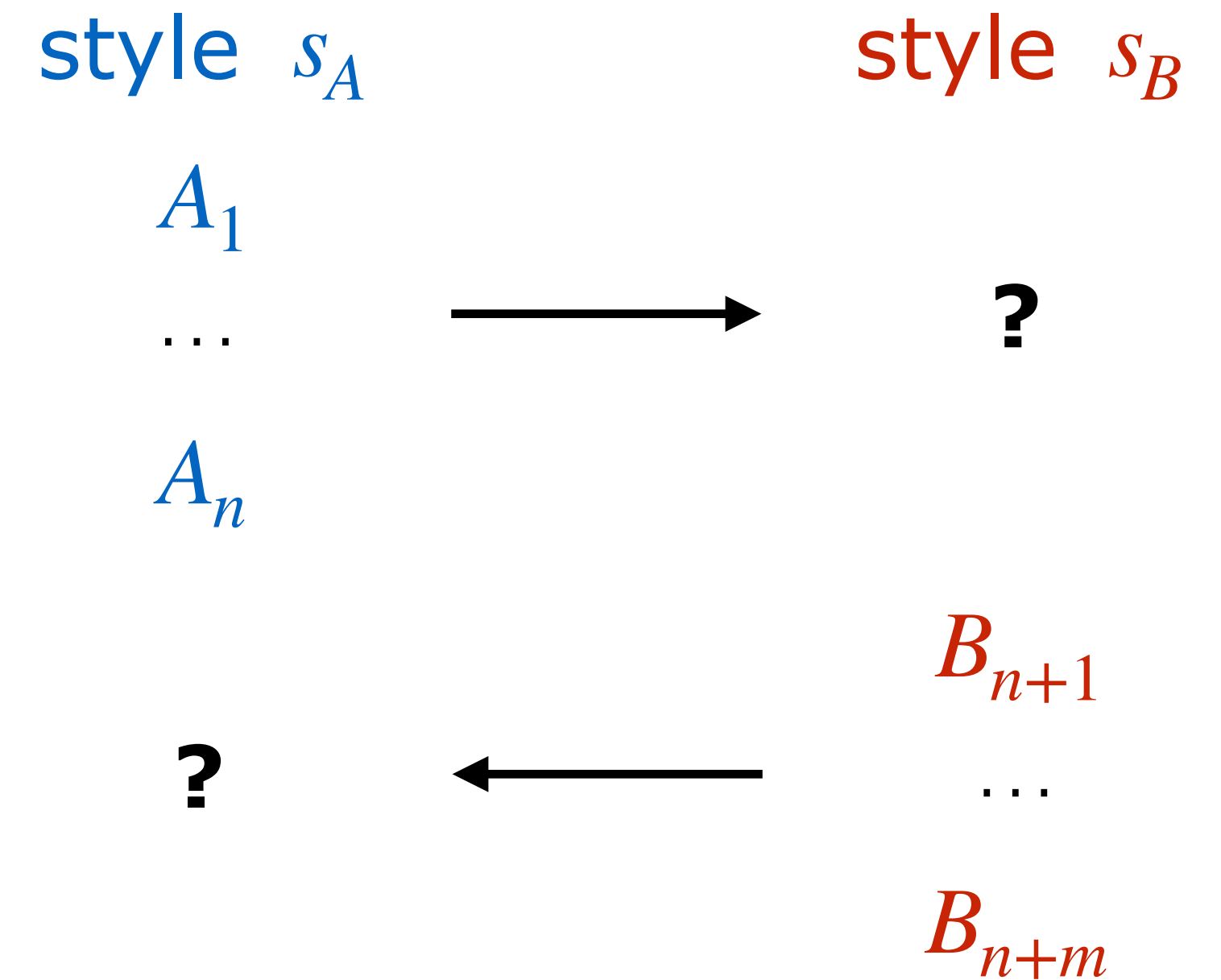
We have three challenges to solve:

1. infer the desired style
2. preserve aspects other than the inferred style
3. generate type of sentences not seen during training

- The task is to transfer style (from the source group to the target positive group), the primary challenge is to transfer style (from the source group to the target positive group), the
- The intra-group variations (category) are **confounding differences** which need to be differentiated from the style and preserved during transfer

Task Formulation

- Given A_1, \dots, A_n of style s_A and B_{n+1}, \dots, B_{n+m} of style s_B , where A_i / B_j is a corpus consisting of sentences x
- Each corpus has its own characteristics
- Change only style and keep other aspects intact



Model Overview

1. Learn a pair of classifiers to detect style and orthogonal attributes
 - Build on invariant risk minimization
2. Use the classifiers to guide a model to transfer in the desired direction

1.0 Invariant Risk Minimization (IRM)

- Specify a set of environments $\mathcal{E} = \{e_1, \dots, e_K\}$, where $e_k = \{(x_k^{(i)}, y_k^{(i)})\}_{i=1}^{n_k}$
- Environment difference accounts for nuisance variation we should **not** pay attention to
- Learn a feature representation that enables the same classifier to be optimal for all environments
- IRMv1: minimize empirical loss across all the data while penalize per-environment gradients with respect to any multiplier of the classifier output

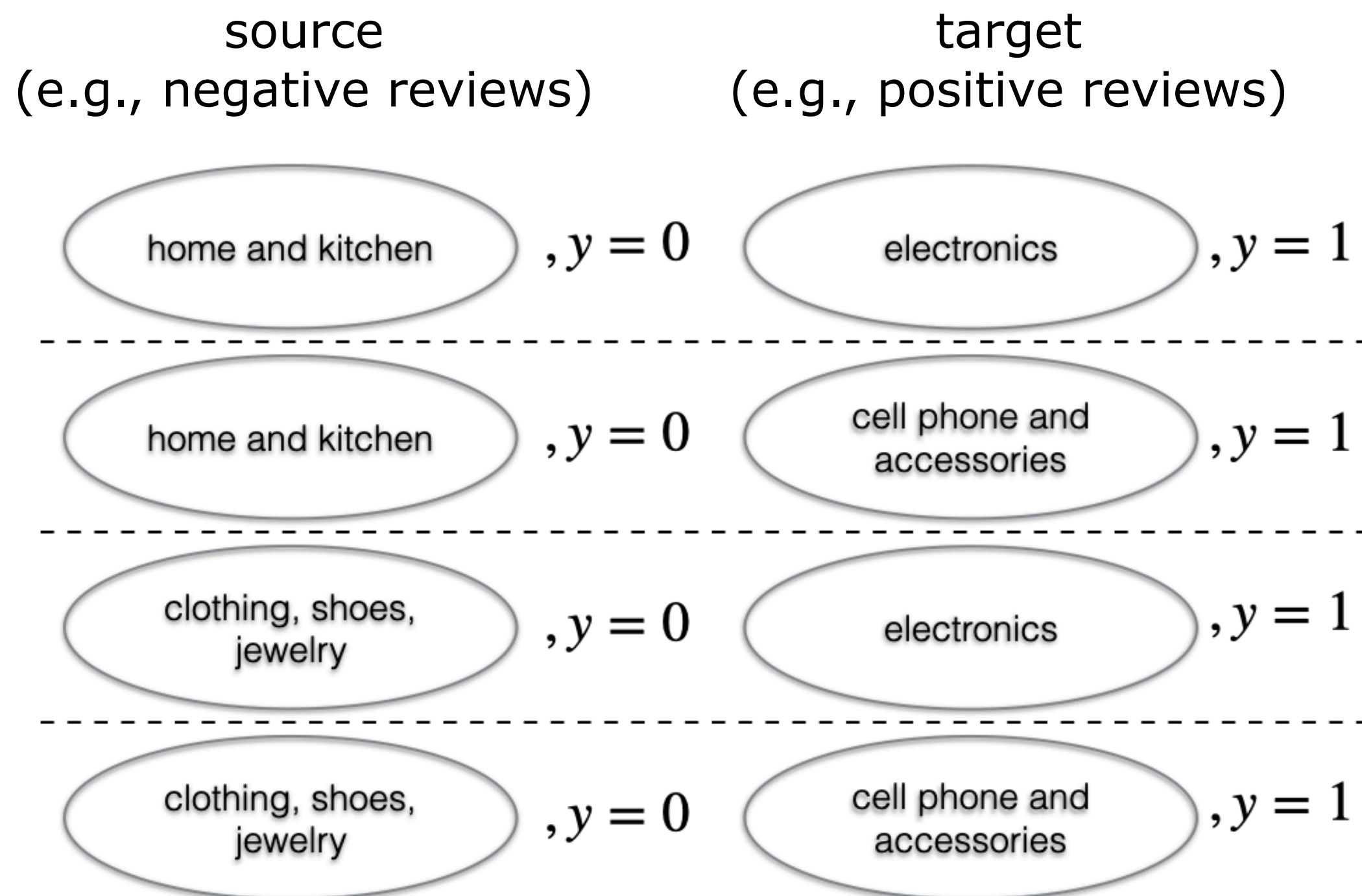
$$\min_{\Phi: \mathcal{X} \rightarrow \mathcal{Y}} \sum_{e \in \mathcal{E}} R^e(\Phi) + \lambda \cdot \|\nabla_{w|w=1.0} R^e(w \cdot \Phi)\|^2$$

$R^e(\Phi) := \mathbb{E}_{X^e, Y^e}[\ell(\Phi(X^e), Y^e)]$
is the risk under environment e

gradients would be zero if Φ
is per-environment optimal

1.1 Inferring Style

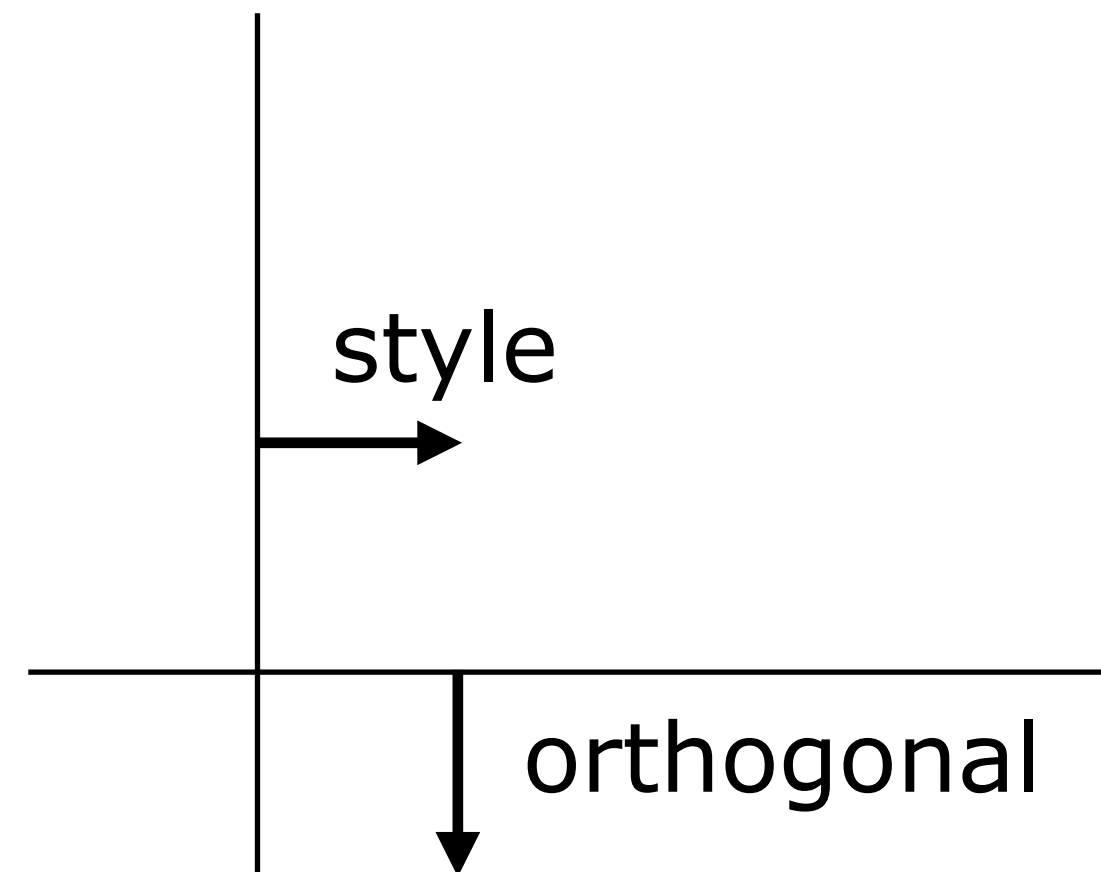
- Construct environments $e_{i,j} = \{(x, y = 0) \mid x \in A_i\} \cup \{(x, y = 1) \mid x \in B_j\}$
- Learn IRM classifier $C_s : \mathcal{X} \rightarrow \mathcal{Y}$ across $\{e_{1,n+1}, \dots, e_{n,n+m}\}$



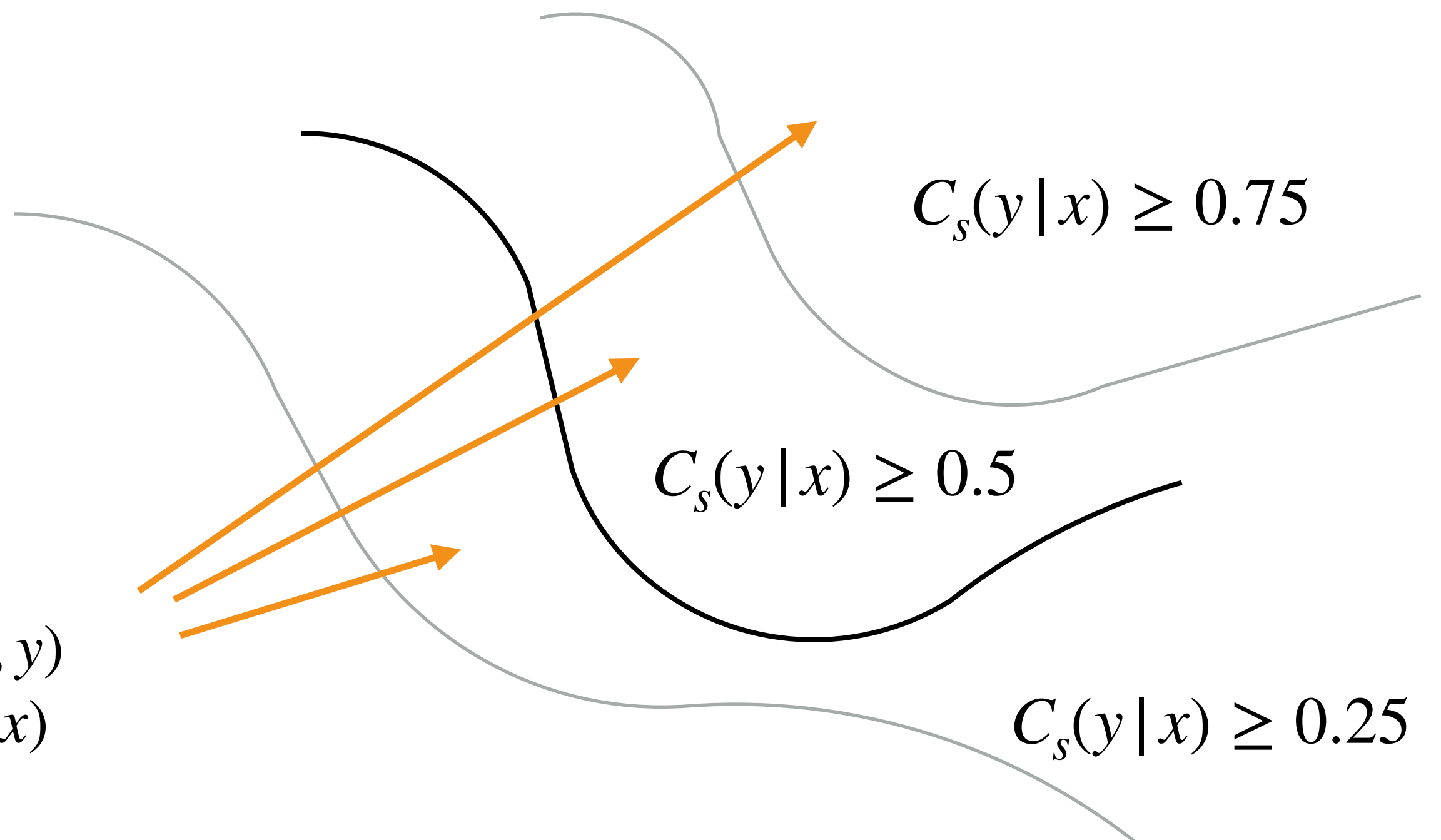
- Since all A_i share style s_A and all B_j share style s_B , style feature elicits an invariant classifier across $\{e_{i,j}\}$
- Conversely, if C_s uses any features specific to A_i/B_j , it won't be optimal in other $e_{i',j'}$

1.2 Inferring Style-Independent Aspects

- Let $A = A_1 \cup \dots \cup A_n$, $B = B_{n+1} \cup \dots \cup B_{n+m}$, $D = \{(x, y = 0) \mid x \in A\} \cup \{(x, y = 1) \mid x \in B\}$
- Construct environments based on C_s :
 $e_1 = \{(x, y) \in D \mid C_s(y|x) > 0.5\}$, $e_2 = \{(x, y) \in D \mid C_s(y|x) \leq 0.5\}$
- Learn IRM classifier $C_o : \mathcal{X} \rightarrow \mathcal{Y}$ across $\{e_1, e_2\}$



$C_o(y|x)$ is invariant
across contours of (x, y)
with respect to $C_s(y|x)$



2. Algorithm for Style Transfer

- Learn $M : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{X}$ that takes a source sentence x and a target group y as input, and outputs a revised sentence that conforms to the style of group y
- Given a data example (x, y) , let $\tilde{x} \sim M(x, 1 - y)$ be the transferred output
 - $\mathcal{L}_{rec} = -\log p_M(x | x, y)$ (reconstruction) → use Gumbel-Softmax to back-propagate
 - temperature annealing
 - length control
 - $\mathcal{L}_{C_s} = -\log p_{C_s}(1 - y | \tilde{x})$ (different style)
 - $\mathcal{L}_{C_o} = -\log p_{C_o}(y | \tilde{x})$ (same orthogonal attributes)
 - $\mathcal{L}_{LM} = \underline{D_{KL}(p_M(\cdot | x, 1 - y) || p_{LM})}$ (language model regularization)
 - $\mathcal{L}_{BT} = -\log p_M(x | \tilde{x}, y)$ maximize entropy (back-translation)

$$\mathbb{E}_{(x,y)}[\mathcal{L}_{rec} + \lambda_1 \mathcal{L}_{C_s} + \lambda_2 \mathcal{L}_{C_o} + \lambda_3 \mathcal{L}_{LM} + \lambda_4 \mathcal{L}_{BT}]$$

Baselines

- M with C_s : without C_o

$$\mathbb{E}_{(x,y)}[\mathcal{L}_{rec} + \lambda_1 \mathcal{L}_{C_s} + \cancel{\lambda_2 \mathcal{L}_{C_o}} + \lambda_3 \mathcal{L}_{LM} + \lambda_4 \mathcal{L}_{BT}]$$

Baselines

- M with C_s : without C_o
- M with C_{ERM} : guided by ERM classifier between A and B instead of C_s and C_o

$$+\lambda \mathcal{L}_{C_{ERM}} = -\log p_{C_{ERM}}(1-y|\tilde{x})$$

$$\mathbb{E}_{(x,y)}[\mathcal{L}_{rec} + \lambda_1 \mathcal{L}_{C_s} + \lambda_2 \mathcal{L}_{C_o} + \lambda_3 \mathcal{L}_{LM} + \lambda_4 \mathcal{L}_{BT}]$$

Baselines

- M with C_s : without C_o
- M with C_{ERM} : guided by ERM classifier between A and B instead of C_s and C_o
- He et al. (2020): regard non-parallel data as partially observed parallel data; treat transferred sequences as latent variables and derive ELBO

$$D_{KL}(p_M(\cdot | x, 1 - y) \| p_{LM_{1-y}})$$

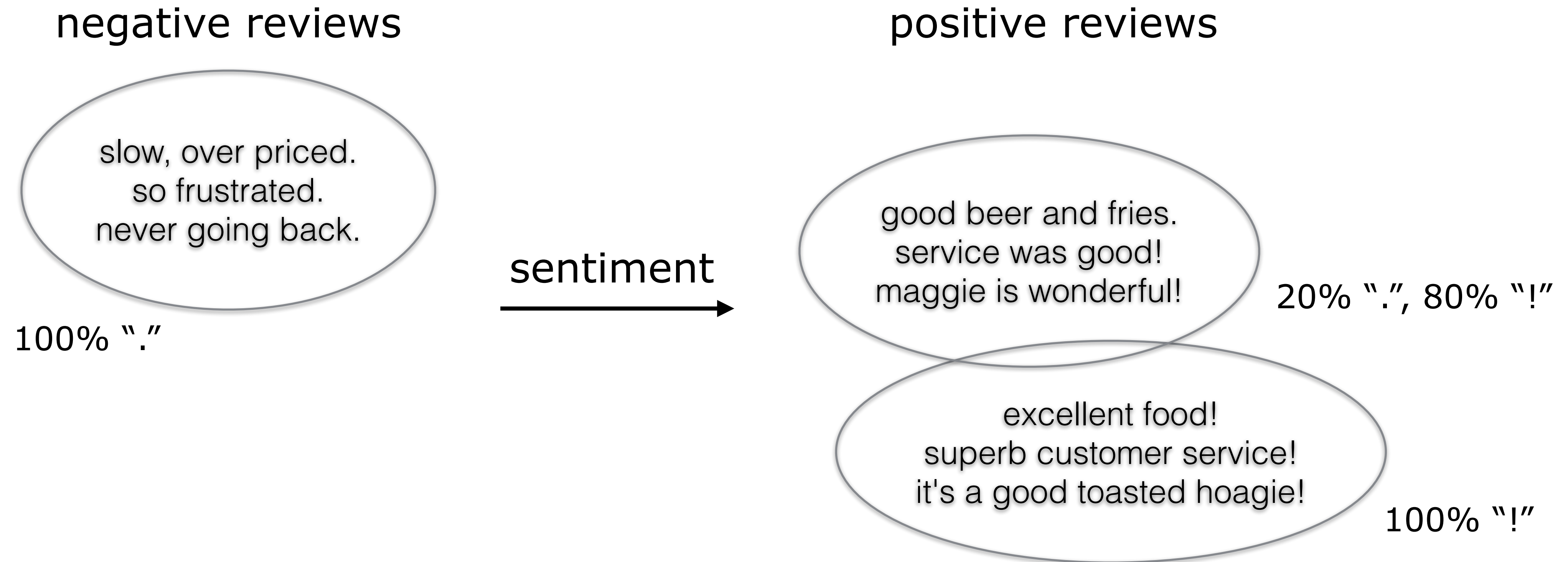
$$\mathbb{E}_{(x,y)}[\mathcal{L}_{rec} + \lambda_1 \mathcal{L}_{C_s} + \lambda_2 \mathcal{L}_{C_o} + \lambda_3 \mathcal{L}_{LM} + \lambda_4 \mathcal{L}_{BT}]$$

Baselines

- M with C_s : without C_o
- M with C_{ERM} : guided by ERM classifier between A and B instead of C_s and C_o
- He et al. (2020): regard non-parallel data as partially observed parallel data; treat transferred sequences as latent variables and derive ELBO
- Krishna et al. (2020): use a separate paraphrasing dataset D_{pp}
 1. train M on D_{pp} , and use it to paraphrase A to A' , B to B'
 2. train inverse models M_A to map A' to A , M_B to map B' to B
 3. to transfer a sentence to style A/B , apply M and then M_A/M_B
 - D_{pp} needs to exclude unwanted changes
 - D_{pp} needs to cover the desired style transformation, otherwise the models are applied OOD

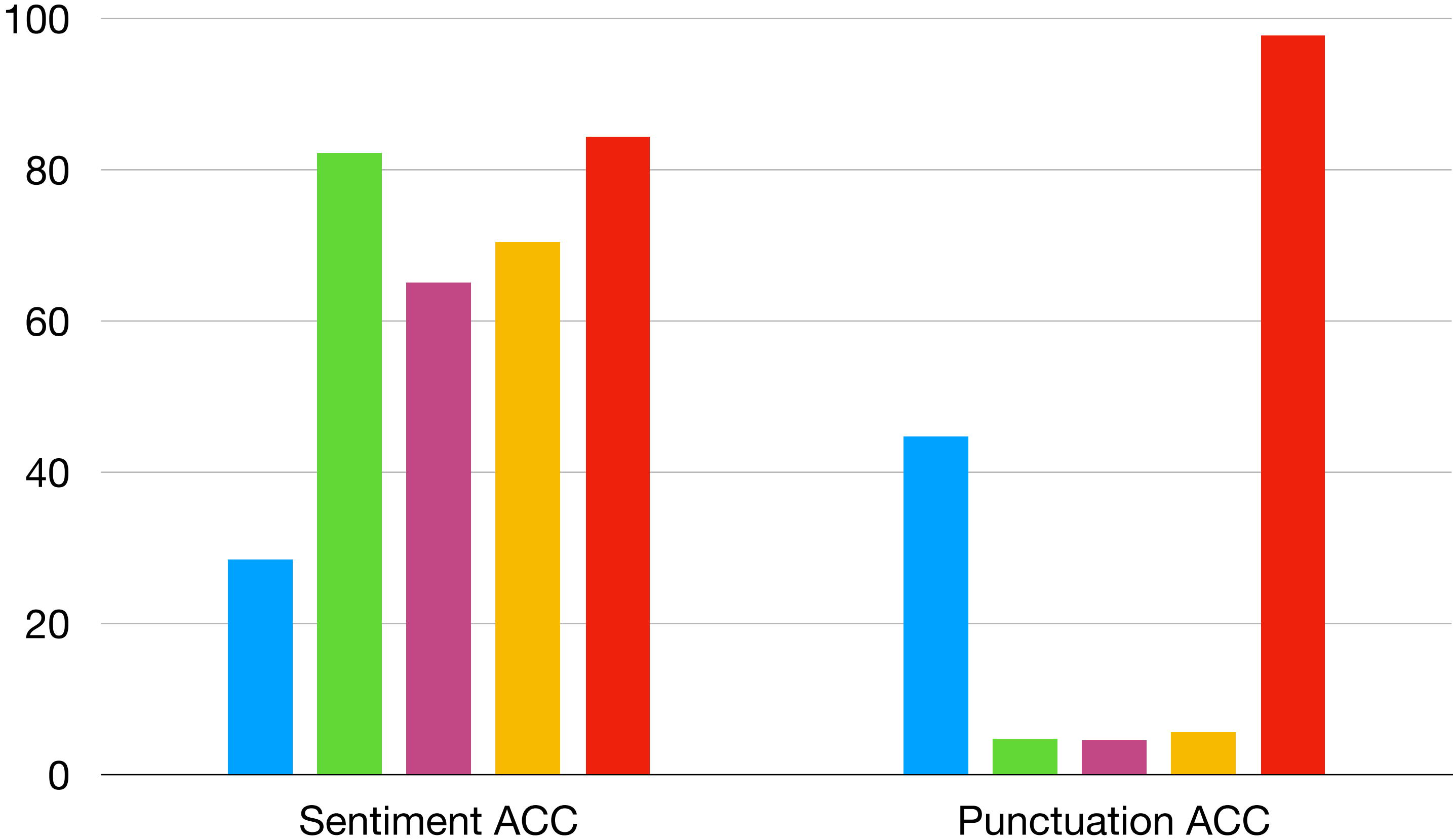
Sentiment Transfer with Different Punctuations

- Adapt sentiment transfer dataset, modifying punctuation to create spurious correlation
- Goal: alter sentiment without changing punctuation



Automatic Evaluation Results

■ Krishna et al. ■ He et al. ■ $M w/ C_{ERM}$ ■ $M w/ C_s$ ■ $M w/ C_s, C_o$



Example Outputs

Input

Krishna et al.

He et al.

Mw/C_{ERM}

Mw/C_s

$Mw/C_s, C_o$ (Ours)

the sales people here are terrible .

the people here are absolutely terrible .

the sales people here are great !

the sales people here are amazing !

the sales people here are fantastic !

the sales people here are amazing .

Input

Krishna et al.

He et al.

Mw/C_{ERM}

Mw/C_s

$Mw/C_s, C_o$ (Ours)

excellent combination of flavors , very unique !

very unique combination of flavors , very unique !” .

horrible customer service .

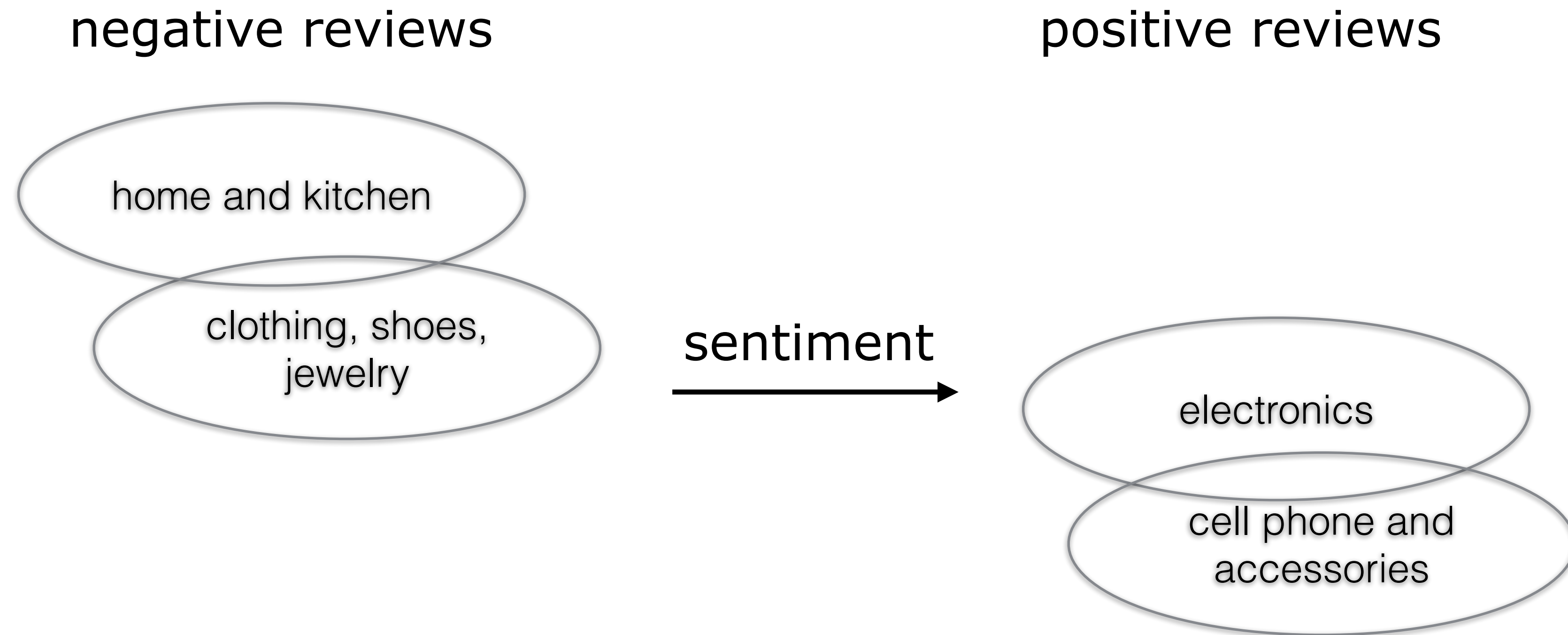
terrible combination of flavors , very disappointing .

terrible combination of flavors , not unique .

terrible combination of flavors , not outstanding !

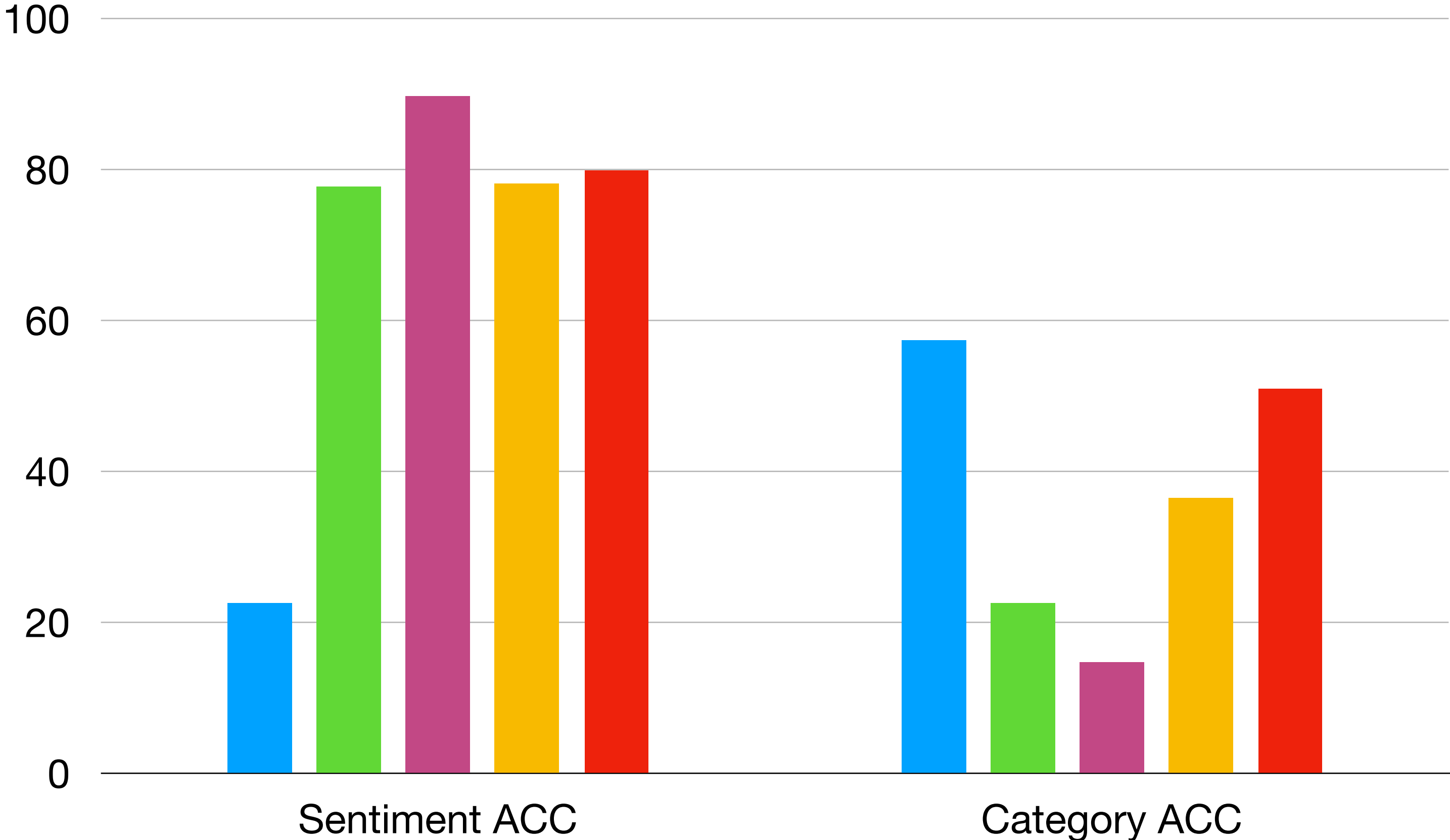
Sentiment Transfer with Different Categories

- Take positive and negative Amazon reviews from different categories
- Goal: alter sentiment without changing product category

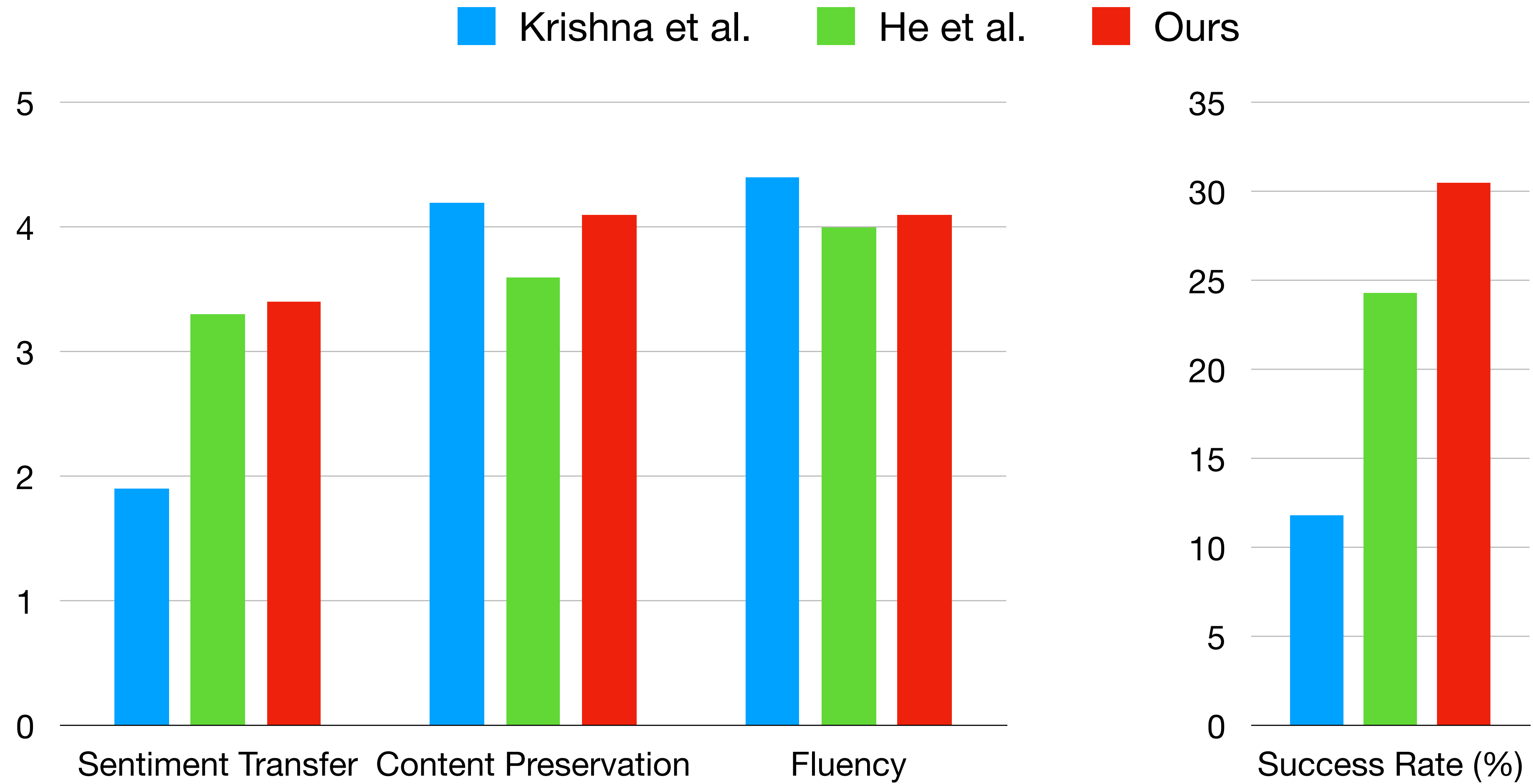


Automatic Evaluation Results

■ Krishna et al. ■ He et al. ■ $M w/ C_{ERM}$ ■ $M w/ C_s$ ■ $M w/ C_s, C_o$



Human Evaluation Results



Example Outputs

Input this shirt was too tight . the sizing seems off .
Krishna et al. the shirt is too tight .
He et al. this case was great . the protection seems great .
Ours this shirt works just perfect . the sizing seems well .

Input the containers do not lock well and are made of low quality materials .
Krishna et al. the containers do not fit securely and are made from poor quality material .
He et al. the phones work well and has made of sound quality of low quality materials .
Ours the containers does the job well and are made of high quality materials .

Input exactly as advertised . converted a molex plug into a sata
Krishna et al. the molex plug was convert to sata as advertised .
He et al. way too big . leaves a inaccurate cut into a bath
Ours not as advertised . converted a molex plug into a sata

A Step Forward: An Aspirational Example

- Transfer from **sonnets** to **tweets** (author is a confounder)

source

target

Shakespeare's **sonnets**
Browning's **sonnets**
Pushkin's **sonnets**
...



? (Shakespeare's **tweets**)

Obama's **tweets**
Bieber's **tweets**
Perry's **tweets**
...