

# Multi-component Similarity Method for Web Product Duplicate Detection

Ronald van Bezu, Sjoerd Borst, Rick Rijkse,  
Jim Verhagen, Damir Vadic, Flavius Frasinca

[vadic@ese.eur.nl](mailto:vadic@ese.eur.nl)

Erasmus University Rotterdam



# Introduction

- E-commerce is growing at a fast pace
- Estimated \$414 billion in 2018 in US
- Entity Resolution (ER)



**SanDisk Sansa Clip+ 8GB Flash MP3 Player**

**\$76.99** from Newegg.com - TOPJ · 1 seller review

SanDisk · Clip+ · 8 GB



**SanDisk Sansa Clip+ 8GB Flash MP3 Player FM Tuner Voice Record...**

**\$109.00** from Fishpond.com · 6 seller reviews

SanDisk · Clip+ · 8 GB · With Radio

# Introduction

- Product descriptions
  - Title
  - Key-Value pairs (basically a `Map[String,String]`)
- Multiple Web shops
- Clean-Clean ER

# Baseline approaches

## TF-IDF approach for ER

- well-known in database ER literature
- Computes the TF-IDF value for each unique term that occurs in the product attribute values
- Each product description is represented by the TF-IDF vector
- Cosine similarity between pairs of vectors -> distance between product descriptions

# Baseline approaches

## Title Model Words Method (TMWM)

- (1) compute cosine similarity between titles;
- (2) if similarity is not high enough, the algorithm extracts *Model Words* from the titles;
- (3) if the non-numeric part is *approximately* the same while the numeric part is not, the two products are classified different;
- (4) else: compute aggregated weighted similarity

# Baseline approaches

- Model words example
  - Samsung - 46" Class/ LED / 1080p / 120Hz / HDTV
  - Samsung - 46" Class/ LED / 1080p / 200Hz / HDTV

# Baseline approaches

## Hybrid Similarity Method (HSM)

- (1) first TMWMM is used
- (2) for matching Key-Value Pairs (KVP), similarity between values is computed
- (3) for non-matching KVP's, model words are extracted and similarity based on % matching is computed
- (4) a final weighted similarity is computed

# Our Approach Similarity

- Multi-component Similarity Method
- All-pair similarity computation
  - scalability addressed in other work based on blocking schemes in a distributed environment



```

1: method KSIM(prod a, prod b)
2:   m := 0
3:   w := 0
4:   I := a.keys and J := b.keys
5:   sim := 0
6:   for all  $k_{i,a} \in a.keys$  do
7:      $k_{i,a} := \text{clean}(k_{i,a})$ 
8:     for all  $k_{j,b} \in b.keys$  do
9:        $k_{j,b} := \text{clean}(k_{j,b})$ 
10:      if keysMatch( $k_{i,a}, k_{j,b}$ ) then
11:         $I := I \setminus k_{i,a}$  and  $J := J \setminus k_{j,b}$ 
12:         $keySim := \text{keySim}(k_{i,a}, k_{j,b})$ 
13:         $valueSim := \text{valSim}(\text{value}(k_{i,a}), \text{value}(k_{j,b}))$ 
14:         $sim := sim + keySim * valueSim$ 
15:         $m := m + 1$ 
16:         $w := w + keySim$ 
17:      end if
18:    end for
19:  end for
20:   $keySim^* := 0$ 
21:  if  $w > 0$  then
22:     $keySim^* := \frac{sim}{w}$ 
23:  end if
24:  return ( $m, I, J, keySim^*$ )
25: end method

```

- ▷ similarity based on keys
- ▷ number of matches
- ▷ weight of matches
- ▷ keys without a match

```

1: method SIM(prod a, prod b)
2:   (m, I, J, keySim*) := KSIM(a, b)
3:   Imw := mw(I) and Jmw := mw(J)
4:   mwSim := mwSim(Imw, Jmw)
5:   titleSim := titleSim(a.title, b.title,  $\alpha$ ,  $\beta$ )
6:   if titleSim = 0 then
7:      $\theta_1 := m / \min(|a.keys|, |b.keys|)$ 
8:      $\theta_2 := 1 - \theta_1$ 
9:   else
10:    
$$\theta_1 := (1 - \mu) \cdot \frac{m}{\min(|a.keys|, |b.keys|)}$$

11:     $\theta_2 := 1 - \mu - \theta_1$ 
12:   end if
13:   sim* :=  $\theta_1 \cdot keySim^* + \theta_2 \cdot mwSim + \mu * titleSim$ 
14:   return sim*
15: end method

```

# Our Approach Clustering

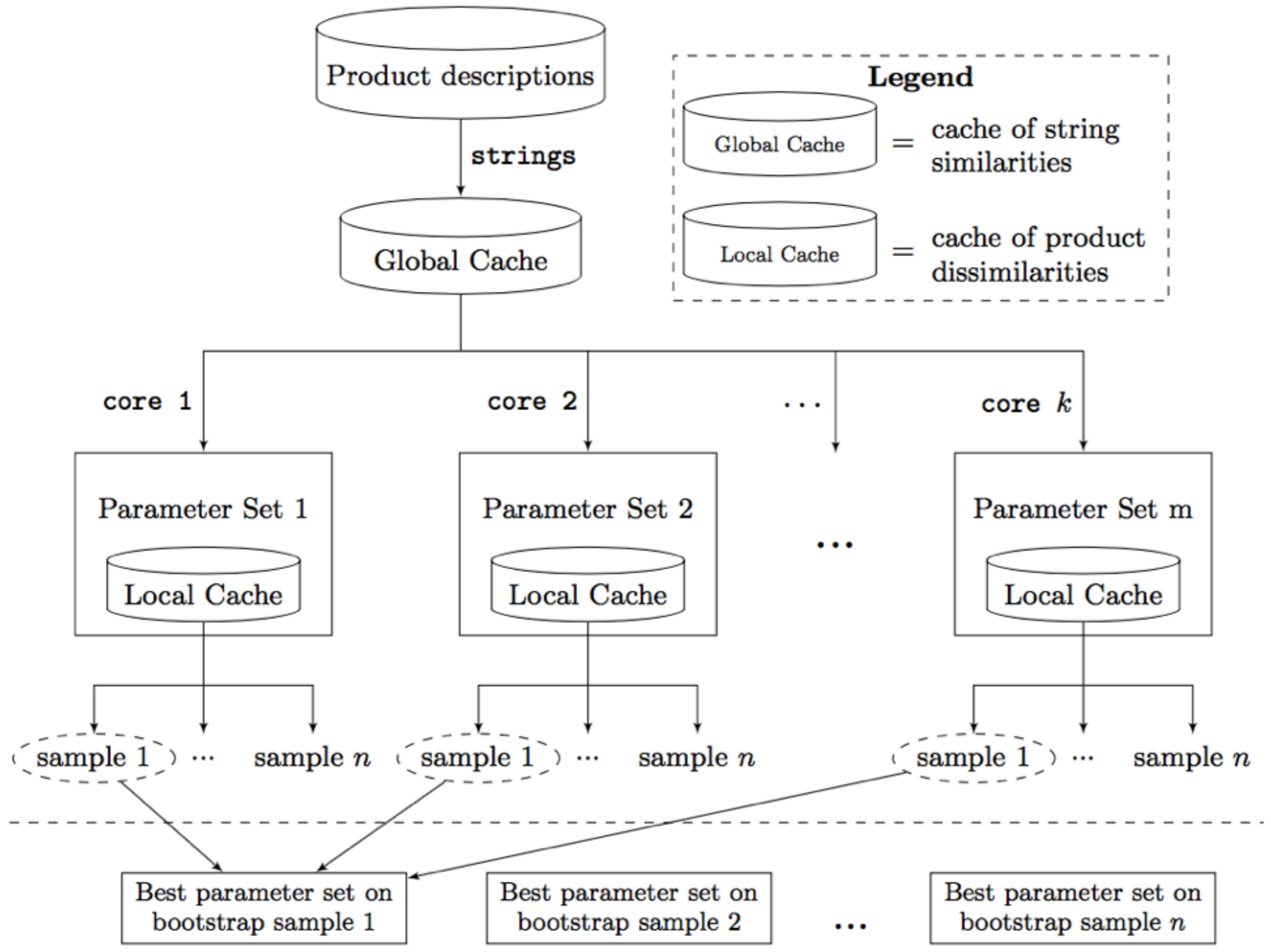
- Hierarchical Clustering
- Handle multiple Web shops at once
- Single linkage with exception to infinite distances to properly handle exclusion cases (e.g., two product descriptions from the same Web shop)

# Evaluation

- Compare TF-IDF, TMWMM, HSM, and MSM
- Data set contains 1629 TV product descriptions
- 1262 of these are unique
- 4 Web shops: Amazon.com, Newegg.com, Best-buy.com, and TheNerds.net
- On average, 29 key-value pairs per product description

# Evaluation

- Interpretation of TP, FP, TN, FN (counting all pairs in produced clusters)
- Baseline algorithms are reimplemented to be used with the hierarchical clustering approach
- Bootstrap samples used as train/test sets
  - trained on FI
- Wilcoxon signed rank test



# Evaluation

- Results

Method	$F_1$ -measure	precision	recall
TF-IDF	0.335	0.337	0.334
TMWWM	0.298	0.349	0.309
HSM	0.287	0.237	0.381
MSM	0.475	0.445	0.512

$p$ -value	TF-IDF	TMWWM	HSM	MSM
TF-IDF	x	1.000	1.000	0.000
TMWWM	0.000	x	0.999	0.000
HSM	0.000	0.001	x	0.000
MSM	1.000	1.000	1.000	x

# Questions?