

# Translating Webpages into Bidphrases for Advertising

Sujith Ravi\*, Andrei Broder†, Evgeniy Gabrilovich†,  
Sandeep Pandey†, Bo Pang†, Vanja Josifovski†

[sravi@isi.edu](mailto:sravi@isi.edu) {[broder](mailto:broder@yahoo-inc.com), [bopang](mailto:bopang@yahoo-inc.com), [gabr](mailto:gabr@yahoo-inc.com), [spandey](mailto:spandey@yahoo-inc.com), [vanjaj](mailto:vanjaj@yahoo-inc.com)}@yahoo-inc.com

\* **University of Southern California**  
**Information Sciences Institute,**  
Los Angeles, CA

† **Yahoo! Research,**  
Santa Clara, CA

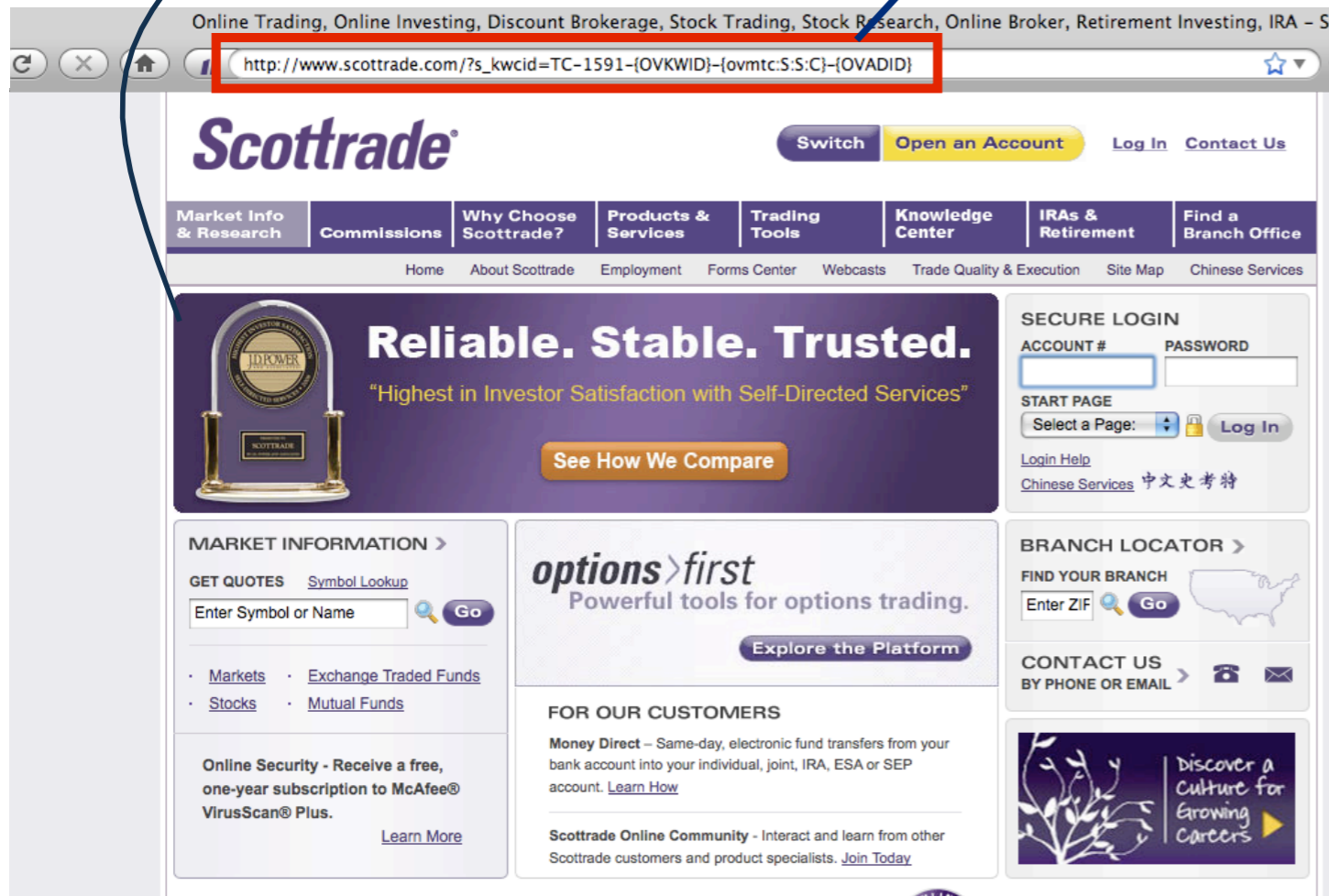
# The Birth of an Ad Campaign

## Landing Page + Landing URL

The image shows a screenshot of a web browser displaying the Scottrade website. The browser's address bar is highlighted with a red box and contains the URL: `http://www.scottrade.com/?s_kwid=TC-1591-{OVKWID}-{ovmtc:S:S:C}-{OVADID}`. Two blue arrows point from the text "Landing Page + Landing URL" above to the address bar and the main content area of the website. The website features the Scottrade logo, navigation links, and a main banner with the text "Reliable. Stable. Trusted." and "Highest in Investor Satisfaction with Self-Directed Services". Other sections include "MARKET INFORMATION", "options>first", "FOR OUR CUSTOMERS", "SECURE LOGIN", "BRANCH LOCATOR", and "CONTACT US".

# The Birth of an Ad Campaign

## Landing Page + Landing URL

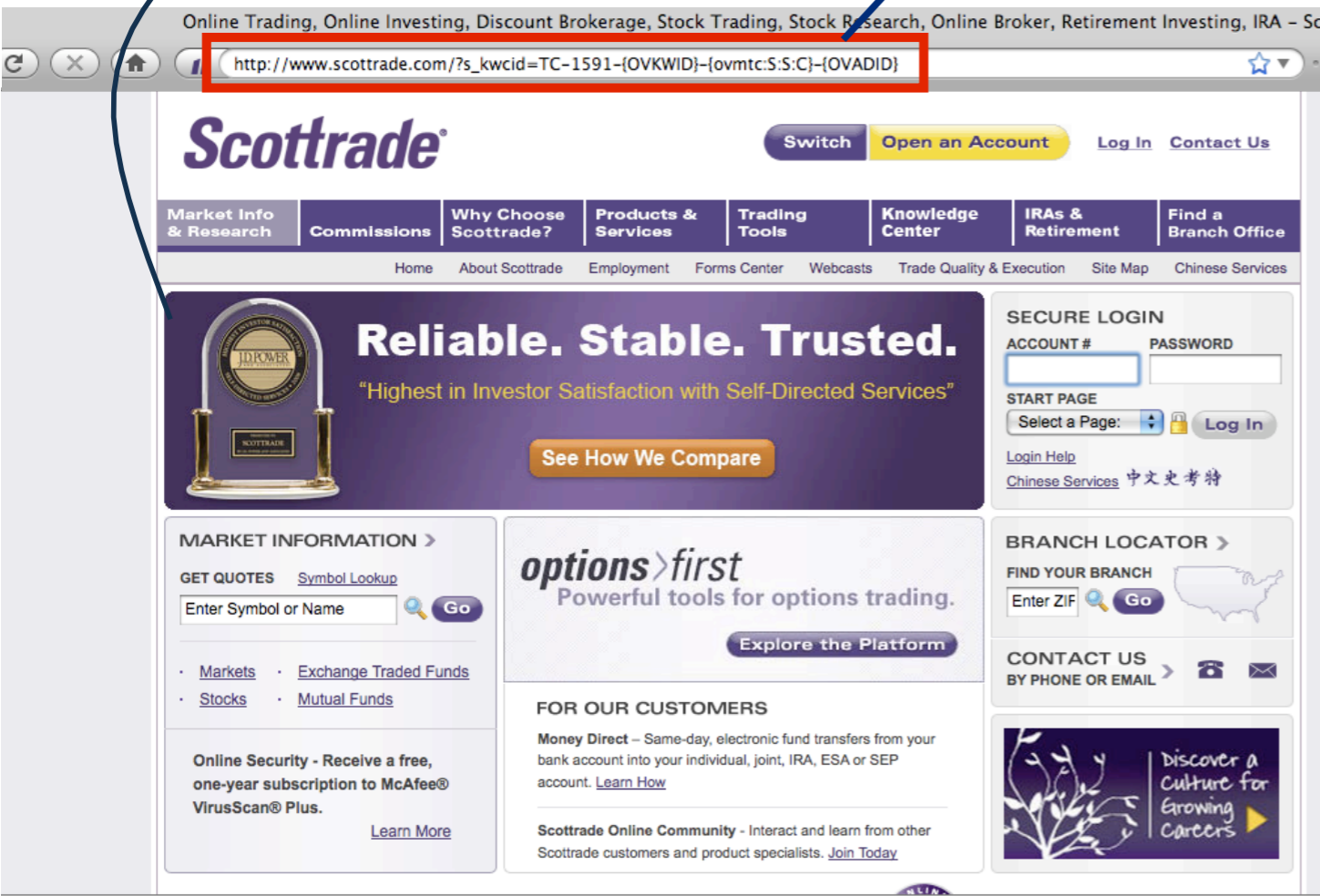


## Relevant Bidphrases

scottrade best brokerage discount  
scottrade switch ira account  
scottrade stock information  
scottrade change ira fund  
scottrade best online brokerage firm  
scottrade best discount broker  
scottrade best online stock broker  
scottrade transfer ira fund  
scottrade switch ira fund

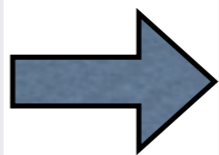
# The Birth of an Ad Campaign

## Landing Page + Landing URL



## Relevant Bidphrases

- scottrade best brokerage discount
- scottrade switch ira account
- scottrade stock information
- scottrade change ira fund
- scottrade best online brokerage firm
- scottrade best discount broker
- scottrade best online stock broker
- scottrade transfer ira fund
- scottrade switch ira fund



Can we automate this process?

# Problem

webpage  
describing a  
product

Landing Page  
 $(\ell)$



Bid Phrases  
 $(b)$

relevant phrases on  
which advertiser can  
potentially bid

# Problem

webpage  
describing a  
product

Landing Page  
( $\ell$ )

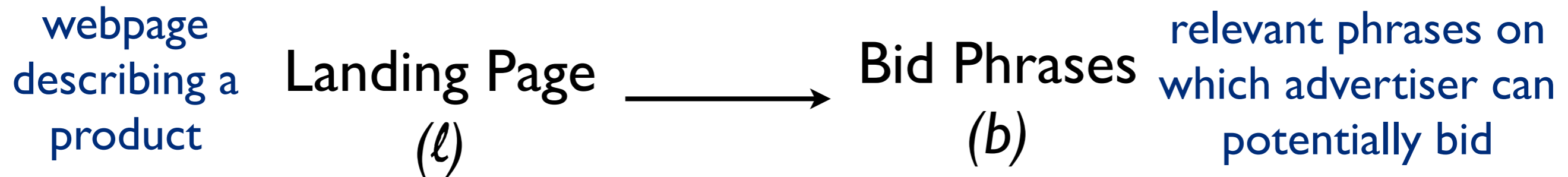


Bid Phrases  
( $b$ )

relevant phrases on  
which advertiser can  
potentially bid

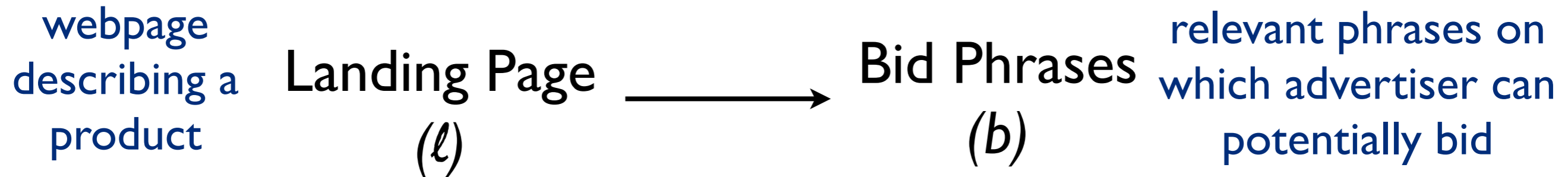
- Can you just get the most informative phrases in the page?

# Problem



- Can you just get the most informative phrases in the page?
- ➔ 96% of ads had at least one bid phrase not in  $\ell$

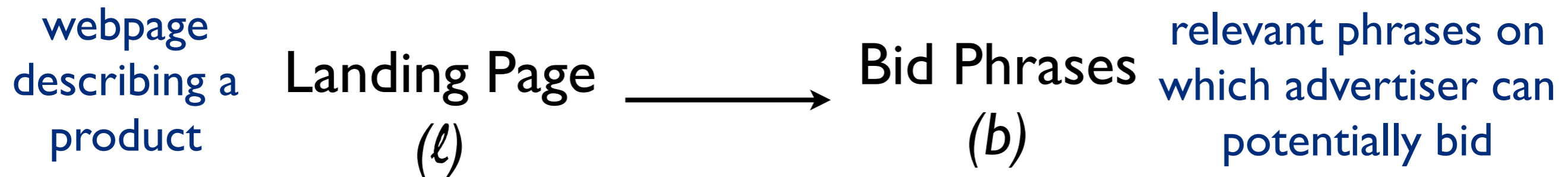
# Problem



- Can you just get the most informative phrases in the page?
- ➔ 96% of ads had at least one bid phrase not in  $\ell$
- How about getting the words?

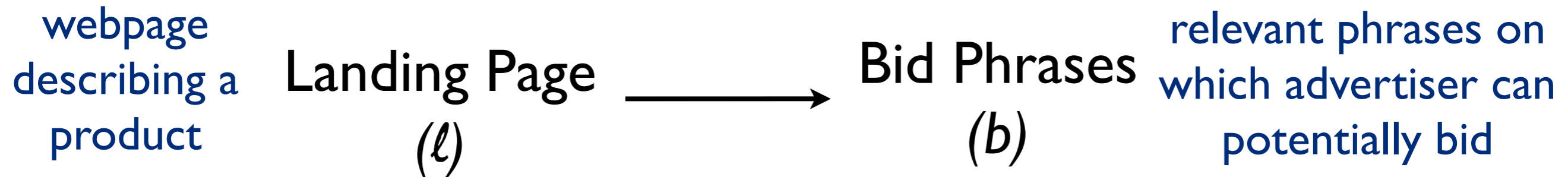


# Problem



- Can you just get the most informative phrases in the page?
- ➔ 96% of ads had at least one bid phrase not in  $\ell$
- How about getting the words?
- ➔ Need to mix-and-match in the right way to generate *phrases*

# Problem



- Can you just get the most informative phrases in the page?
- ➔ 96% of ads had at least one bid phrase not in  $\ell$
- How about getting the words?
- ➔ Need to mix-and-match in the right way to generate *phrases*
- ➔ The bid phrase set for 70% of ads contained one or more words not in  $\ell$

# A Two-phase Approach

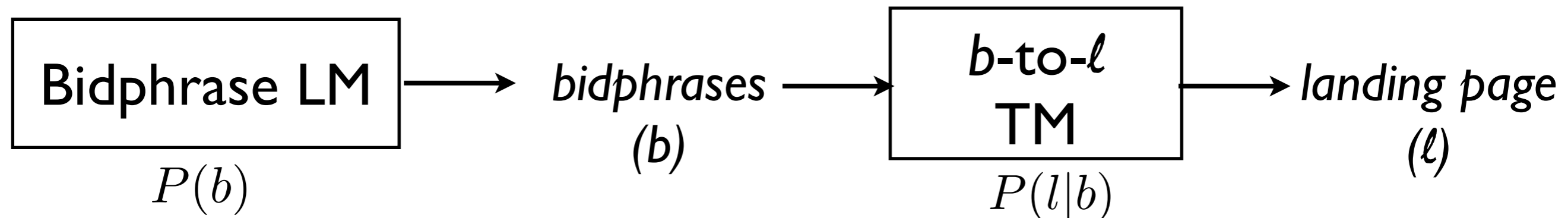
1. Candidate bid phrases are generated
  - need to be able to generate “novel” phrases
2. Candidates are ranked
  - need to pick phrases relevant to page and resemble queries

# Translation-based Approach

Landing Page + Landing URL  $(l)$   $\xrightarrow{\text{translation}}$  bidphrases  $(b)$

- Noisy-channel approach used in Machine Translation

## Generative Model

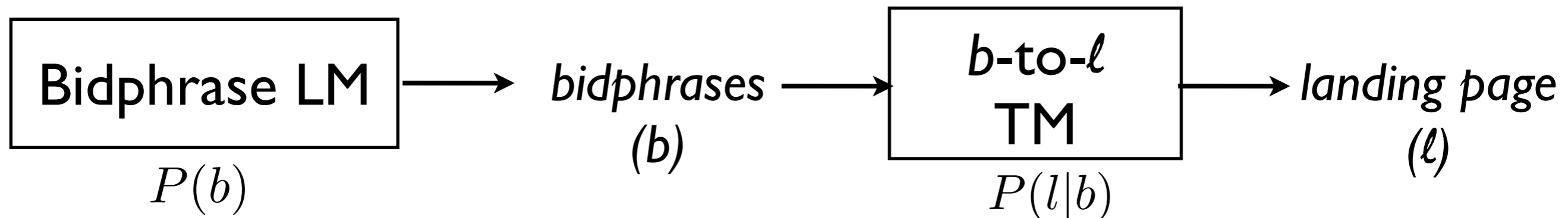


# Translation-based Approach

Landing Page + Landing URL  $(l)$   $\xrightarrow{\text{translation}}$  bidphrases  $(b)$

- Noisy-channel approach used in Machine Translation

## Generative Model



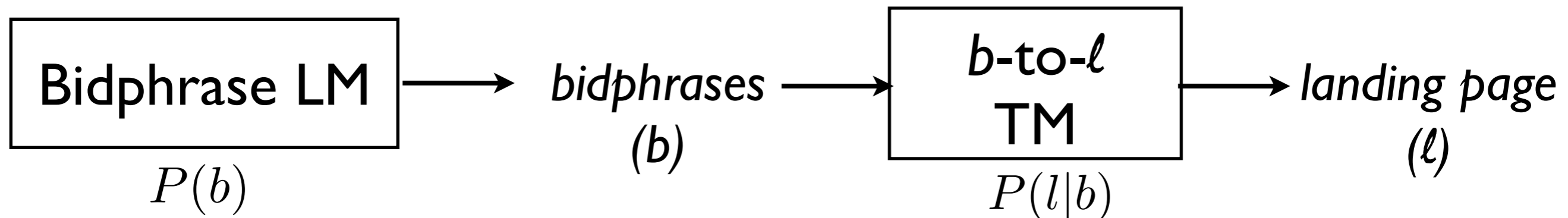
Language Model generates potential bidphrases

# Translation-based Approach

Landing Page + Landing URL  $(l)$   $\xrightarrow{\text{translation}}$  bidphrases  $(b)$

- Noisy-channel approach used in Machine Translation

## Generative Model



Language Model generates potential bidphrases

Translation Model translates each bidphrase word  $(b_i)$  into a word  $(l_i)$  appearing on the landing page

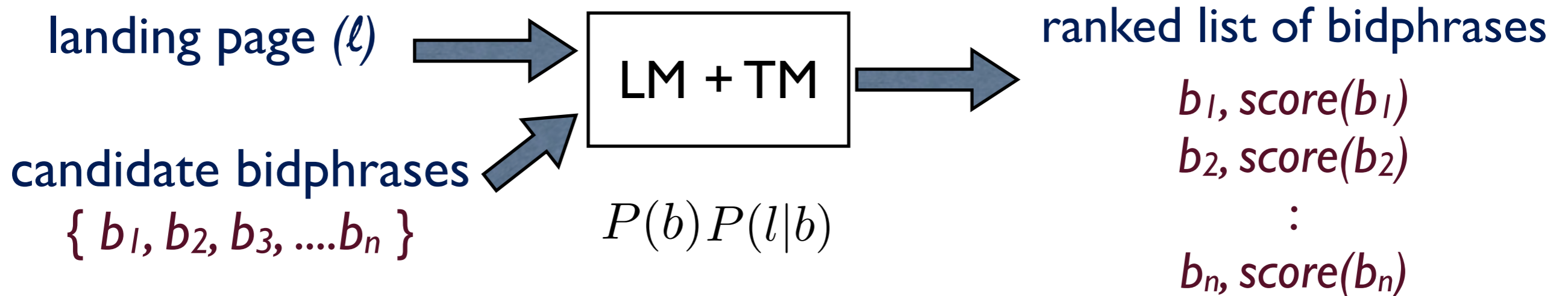
# Ranking Candidate Phrases

Given candidate bidphrases +

Bidphrase LM

*b-to-l*  
TM

Score candidate bidphrases (Decoding)



# Bidphrase Language Model LM

Bidphrases should resemble queries

Estimating the model

- LM is a bigram language model, with *back-off* to a unigram model
- Model estimated on a large query corpus  $Q$  (*~76 million queries from Yahoo! Web search log*)



# Translation Model

TM

Estimating the model

IBM Model I

- Estimate translation table  $t(l_j | b_i)$  to maximize likelihood of (parallel) data (bid phrase, page) pairs

$$\Pr(l|b) \propto \prod_j \sum_i t(l_j | b_i)$$

$l_j$  = word in landing page  $l$

$b_i$  = word in bidphrase  $b$

# Translation Model

TM

Estimating the model

IBM Model I

- Estimate translation table  $t(l_j | b_i)$  to maximize likelihood of (parallel) data (bid phrase, page) pairs

$$\Pr(l|b) \propto \prod_j \sum_i t(l_j | b_i)$$

$l_j$  = word in landing page  $l$

$b_i$  = word in bidphrase  $b$

- *Null* token added to bidphrase side to account for irrelevant words from landing page

# Translation Model

TM

Estimating the model

IBM Model I

- Estimate translation table  $t(l_j | b_i)$  to maximize likelihood of (parallel) data (bid phrase, page) pairs

$$\Pr(l|b) \propto \prod_j \sum_i t(l_j | b_i)$$

$l_j$  = word in landing page  $l$

$b_i$  = word in bidphrase  $b$

- *Null* token added to bidphrase side to account for irrelevant words from landing page
- Incorporate importance of words in a page

$$\Pr(l|b) \propto \prod_j (\sum_i t(l_j | b_i))^{w_j}$$

$w_j$  = importance weight assigned to word  $l_j$

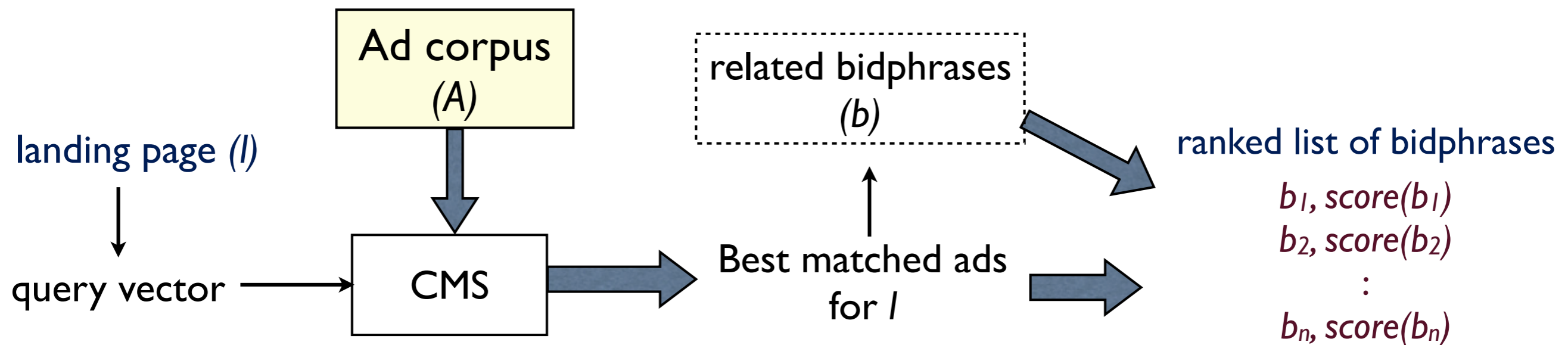
(higher weight for words appearing in titles, headings, etc.)

# Generating Candidate Phrases

- Theoretically, all phrases in query log can be candidates
  - ➔ inefficient
- Strategy-1: Build a candidate set containing only phrases appearing on landing page  $\{b_{LP}\}$ 
  - ➔ downside: no novel bidphrases generated
- Strategy-2: Use translation model (TM) to generate novel bidphrases  $\{b_{TMgen}\}$ 
  - ➔ bridges vocabulary mismatch
  - ➔ use only *salient* words from landing page to generate new candidates

# Alternative Methods

- Extraction-based system (Baseline)
  - extract candidates from page, rank by  $\text{cosine-sim}(b, \ell)$
- Discriminative system using  $\text{SVM}^{\text{rank}}$  using features:
  - *word-overlap, position on page, cosine-sim(b, l), ...*
- Using content-match system (CMS)



# Evaluation

## Large-scale and automatic?

For each landing page ( $\ell$ ) in the test corpus (10,500 pages)

- Gold-standard bidphrases  $\{b_{gold}\}$ 
  - provided by the advertisers
  - average 9 per landing page
- Each generated bidphrase ( $b_c$ ) is compared against  $\{b_{gold}\}$

Relative ordering should be meaningful

# Evaluation Metrics

$b_c$  against  $\{b_{gold}\}$  of page  $l$

## 1. Minimum Edit Distance (*minED*)

$$\text{minED}(b_c, l) = \min_{b_j \in \{b_{gold}\}} \text{ED}(b_c, b_j)$$

where,

$$\text{ED}(b_c, b_j) = \frac{\# \text{ of oprns. to convert } b_c \rightarrow b_j}{\# \text{ of words in } b_j}$$

lower *minED* scores => better

## 2. ROUGE-1 metric

$$\text{ROUGE-1}(b_c, l) = \frac{\sum_{b_j \in \{b_{gold}\}} \# \text{ of words in } b_c \cap b_j}{\sum_{b_j \in \{b_{gold}\}} \# \text{ of words in } b_j}$$

higher *ROUGE-1* scores => better

# Evaluation Metrics

$b_c$  against  $\{b_{gold}\}$  of page  $l$

## 1. Minimum Edit Distance (*minED*)

$$\text{minED}(b_c, l) = \min_{b_j \in \{b_{gold}\}} \text{ED}(b_c, b_j)$$

where,

$$\text{ED}(b_c, b_j) = \frac{\# \text{ of oprns. to convert } b_c \rightarrow b_j}{\# \text{ of words in } b_j}$$

lower *minED* scores => better

## 2. ROUGE-1 metric

$$\text{ROUGE-1}(b_c, l) = \frac{\sum_{b_j \in \{b_{gold}\}} \# \text{ of words in } b_c \cap b_j}{\sum_{b_j \in \{b_{gold}\}} \# \text{ of words in } b_j}$$

higher *ROUGE-1* scores => better

Is it similar to any phrase in  $\{b_{gold}\}$  ?



# Evaluation Metrics

$b_c$  against  $\{b_{gold}\}$  of page  $l$

## I. Minimum Edit Distance (*minED*)

$$\text{minED}(b_c, l) = \min_{b_j \in \{b_{gold}\}} \text{ED}(b_c, b_j)$$

where,

$$\text{ED}(b_c, b_j) = \frac{\# \text{ of oprns. to convert } b_c \rightarrow b_j}{\# \text{ of words in } b_j}$$

lower *minED* scores => better

Is it similar to any phrase in  $\{b_{gold}\}$  ?

## 2. ROUGE-1 metric

$$\text{ROUGE-1}(b_c, l) = \frac{\sum_{b_j \in \{b_{gold}\}} \# \text{ of words in } b_j}{\sum_{b_j \in \{b_{gold}\}} \# \text{ of words in } b_j}$$

recall of words in  $\{b_{gold}\}$

higher *ROUGE-1* scores => better

# Main Comparisons



$\{b_{LP}\}$

$\{b_{LP+CMS}\}$

$\{b_{LP+TM_{gen}}\}$

Candidate  
Generation

# Main Comparisons



$\{b_{LP}\} \rightarrow$  words/phrases extracted from landing page

$\{b_{LP+CMS}\}$

$\{b_{LP+TM_{gen}}\}$

Candidate  
Generation

# Main Comparisons



$\{b_{LP}\} \Rightarrow$  words/phrases extracted from landing page

$\{b_{LP+CMS}\} \Rightarrow$  + bidphrases proposed by CMS

$\{b_{LP+TM_{gen}}\}$

Candidate  
Generation

# Main Comparisons



$\{b_{LP}\} \Rightarrow$  words/phrases extracted from landing page

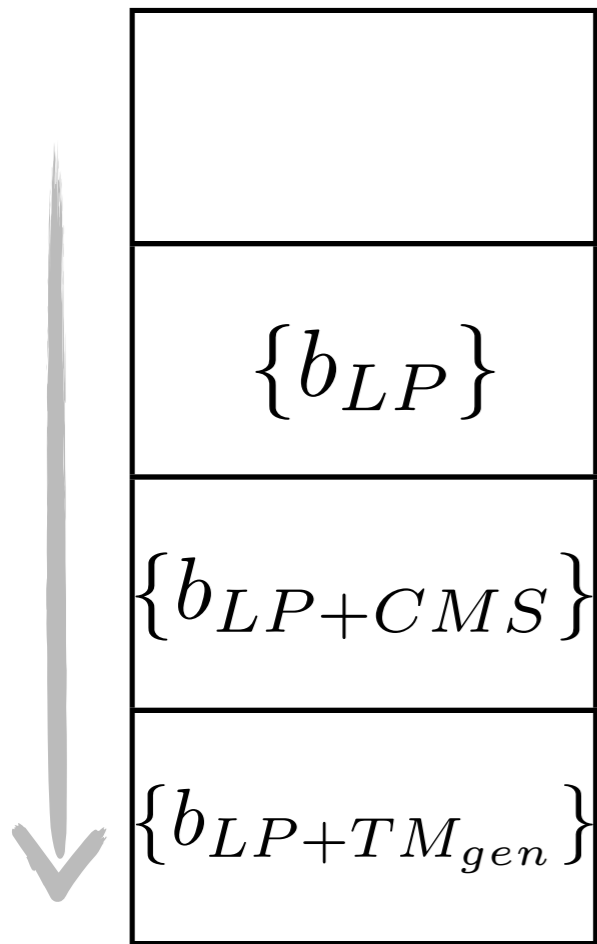
$\{b_{LP+CMS}\}$

$\{b_{LP+TM_{gen}}\} \Rightarrow$  + new phrases generated by translating landing page content using TM

Candidate  
Generation

# Main Comparisons

Candidate Ranking



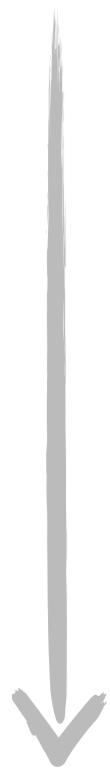
Candidate  
Generation

# Main Comparisons

Candidate Ranking



	cosine
$\{b_{LP}\}$	baseline
$\{b_{LP+CMS}\}$	
$\{b_{LP+TM_{gen}}\}$	



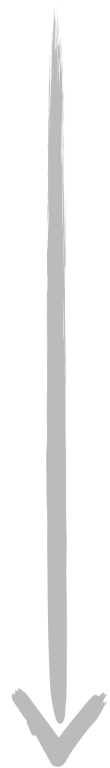
Candidate  
Generation

# Main Comparisons

Candidate Ranking



	cosine	CMS
$\{b_{LP}\}$	baseline	
$\{b_{LP+CMS}\}$		CMS
$\{b_{LP+TM_{gen}}\}$		



Candidate  
Generation

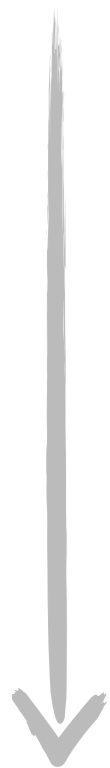


# Main Comparisons

Candidate Ranking



	cosine	CMS	$SVM^{rank}$
$\{b_{LP}\}$	baseline		
$\{b_{LP+CMS}\}$		CMS	Discriminative system
$\{b_{LP+TM_{gen}}\}$			



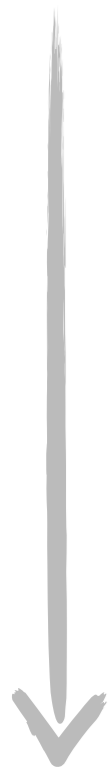
Candidate  
Generation

# Main Comparisons

Candidate Ranking



	cosine	CMS	$SVM^{rank}$	LM+TM
$\{b_{LP}\}$	baseline			
$\{b_{LP+CMS}\}$		CMS	Discriminative system	
$\{b_{LP+TM_{gen}}\}$				Translation-based



Candidate  
Generation

# Main Comparisons

lower *minED* scores => better bidphrases

	Baseline (cosine)	CMS	Discriminative System (SVM <sup>rank</sup> with features)	LM+TM <i>B<sub>LP+TM<sub>gen</sub></sub></i>
minED @ rank 1	0.66	0.78	0.67	0.68
minED @ rank 5	0.71	0.81	0.72	0.68
minED @ rank 10	0.75	0.83	0.74	0.70
ROUGE-1 @ rank 1	0.24	0.22	0.26	0.29
ROUGE-1 @ rank 5	0.19	0.21	0.24	0.28
ROUGE-1 @ rank 10	0.16	0.20	0.22	0.27

higher *ROUGE-1* scores => better bidphrases

Test corpus = 10,500 landing pages

# Main Comparisons

	cosine	CMS	SVM <sup>rank</sup>	LM+TM
$\{b_{LP}\}$				
$\{b_{LP+CMS}\}$				
$\{b_{LP+TM_{gen}}\}$				

Test corpus = 10,500 landing pages

# Candidate Generation

	cosine	CMS	SVM <sup>rank</sup>	LM+TM
$\{b_{LP}\}$				
$\{b_{LP+CMS}\}$				
$\{b_{LP+TM_{gen}}\}$				

Test corpus = 10,500 landing pages

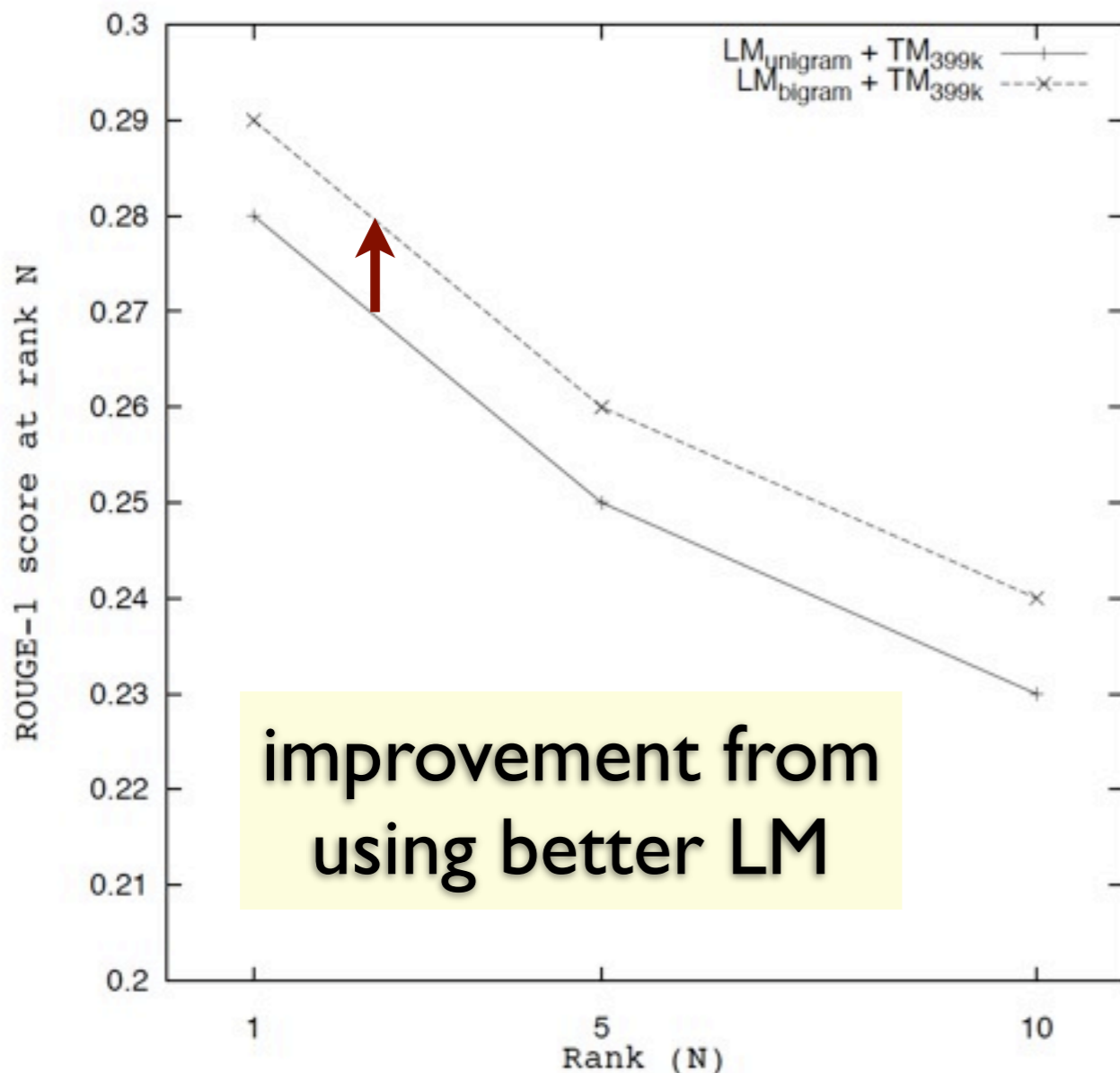
# Ranking Methods

	cosine	CMS	$SVM^{rank}$	LM+TM
$\{b_{LP}\}$				
$\{b_{LP+CMS}\}$				
$\{b_{LP+TM_{gen}}\}$				

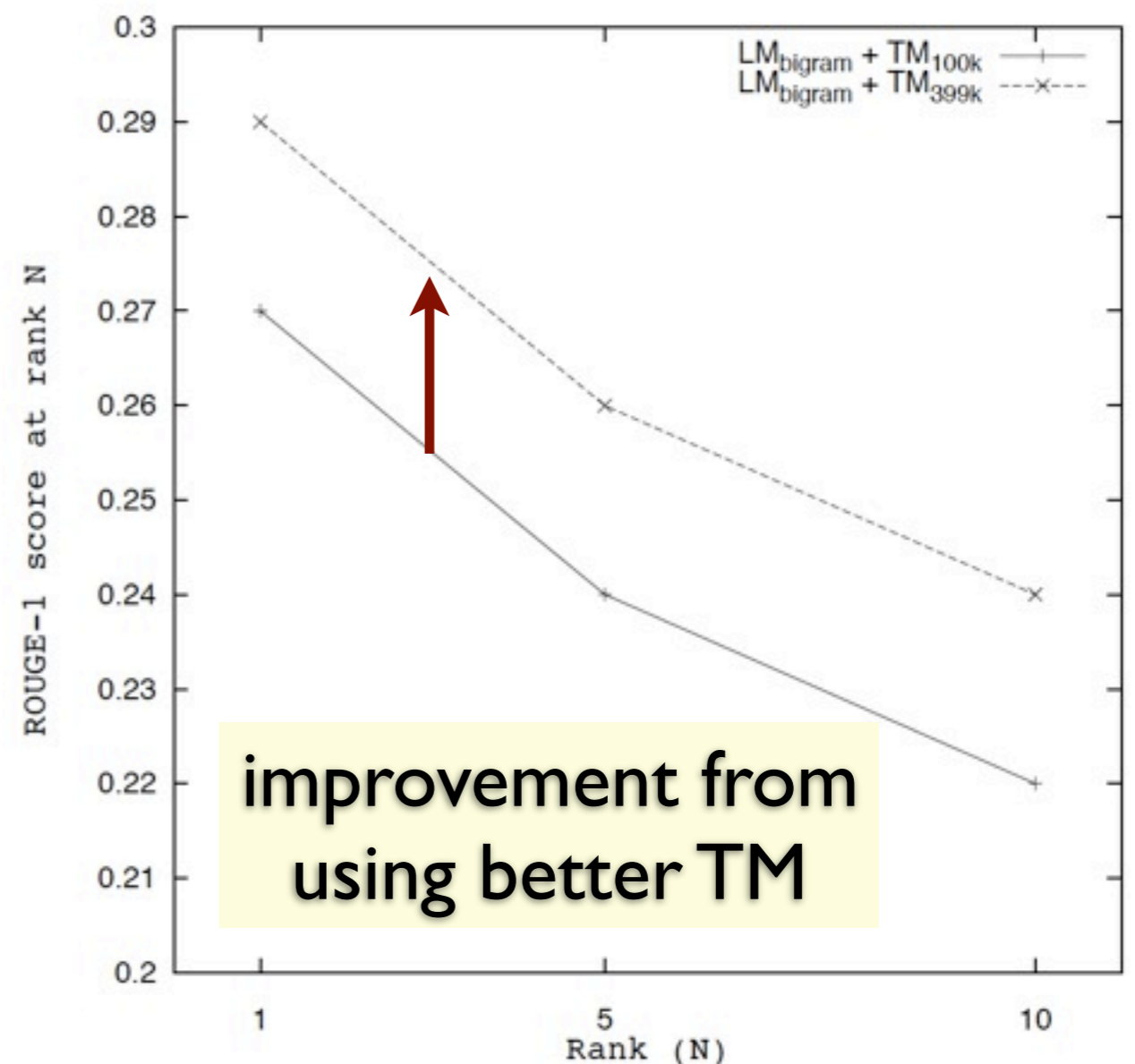
Test corpus = 10,500 landing pages

# Component Analysis for TM+LM

Varying the Language Model  
(unigram vs. bigram LM)



Varying the Translation Model  
(using different training sizes)



# Translation Table

bidphrase word ( $b_i$ )	top translations	
	$l_j$	$P(l_i b_j)$
account	<i>account</i>	0.756
	<i>accounts</i>	0.023
	<i>checking</i>	0.012
	<i>savings</i>	0.010
	<i>online</i>	0.009
addiction	<i>addiction</i>	0.940
	<i>drug</i>	0.012
	<i>alcohol</i>	0.009
	<i>rehab</i>	0.007
	<i>addict</i>	0.005

bidphrase word ( $b_i$ )	top translations	
	$l_j$	$P(l_i b_j)$
mag	<i>mag</i>	0.870
	<i>magazine</i>	0.025
	<i>cover</i>	0.013
	<i>subscription</i>	0.009
	<i>magazin</i>	0.008
ticket	<i>ticket</i>	0.288
	<i>tickets</i>	0.220
	<i>flights</i>	0.037
	<i>prices</i>	0.030
	<i>fares</i>	0.020



# Related Work

- Online Advertising
  - keyword extraction [Yih et al., 2006]
  - bridging vocabulary overlap in contextual advertising [Ribeiro-Neto et al., 2005]
  - query expansion and rewriting, keyword suggestion, ...
- Machine Translation / noisy channel model
  - text summarization [Knight and Marcu, 2000];  
paraphrase extraction [Quirk et al., 2004]
  - contextual advertising [Murdock et al., 2007]

# Conclusion

- Several automatic methods to generate bidphrases for online advertising
- Two evaluation measures proposed to assess different qualitative aspects of generated bidphrases
- **Novel translation-based approach using a generative model**
  - ➔ produces best results in terms of both evaluation measures
  - ➔ generates novel phrases that are relevant but do not appear on page