# Pre-computing Search Features for Fast and Accurate Query Classification

**Arnd Christian König**

**Venkatesh Ganti**

**Xiao Li**

**Microsoft Research**

# Query Classification

# Query Classification



**Query Classification important for**
- ■ Matching advertisements for a query
- ■ Retrieval of additional non-web search results from verticals
- ■ Load optimization for search verticals

# Problem Statement

**Query Classification:**

- **Problem**: ~3 words in queries => little 'signal' for classification.
- Large vocabulary size => large, sparse feature space.
- Difficult to generalize across queries.

**Post-Retrieval Features:**

- Use search to obtain more context to derive features.

# Post-Retrieval features

**General approach:**

- Issue the search query against a document corpus .
- Identify relevant sub-components of top results (e.g., titles, captions, key terms, etc.)
- Derive additional features from these components.



$$F = \sum_{d \in \text{Result}} (\cdots)$$

# Problem Statement

**Query Classification:**

- **Problem:** ~3 words in queries => little 'signal' for classification.
- Large vocabulary size => large, sparse feature space.
- Difficult to generalize across queries.

**Post-Retrieval Features:**

- Use search to obtain more context to derive features.
- => significant improvements in classification accuracy.

- **Problem:  Search Latency**
  - Even slight (100 ms) increases in latency decrease user satisfaction, increase in fraction of abandoned searches.

$\Rightarrow$ **Task:** Realize benefits of post-retrieval features at low overhead.

# Our Approach

**Query: Low-light snapshots**

**Classification Task: Product-intent**
**Tags: Entity Categories**

**Features based on the incidence of tags in the documents returned in response to a query.**

$\Rightarrow$ **Small feature space, features generalize across queries.**
$\Rightarrow$ **Less information to store, helping pre-computation.**

**Tag-Ratios**

$q, t) :=$

$t \in result(q)$

$esult(q)|$

**Other examples:**
*Corpus:* **Sponsored Search Bids**
- *Tags:* **Advertiser-IDs**
- **Advertisers can be thought of as 'topics'**

*Corpus:* **Wikipedia**
- *Tags:* **Wikipedia-Category Tags**

# Our Approach (II)

Documents $\mathcal{C}$

Tag Corpus $\mathcal{T}$

**Pre-computation of Tag-ratios**

**Retrieval Semantics: word-containment**
- Search engine not involved in retrieval
  $\Rightarrow$ Fast pre-computation of query sets
- Tradeoff: result relevance vs. result size

Collection of (query, tag-ratio) pairs

e-computed and indexed in memory

**Feature Generation**

**Query Classifier**

Online

**The rest of this talk:**
- How do we generate features from the ratios?
- Size-constraints: for which queries do we pre-compute ratios?
- How do we deal with query that is not pre-computed?

# Creating Features from Tag Ratios

**Features = Ratios?**

$$F = \left[ratio(q, t_1), \cdots, ratio(q, t_k)\right]$$

- **Problem I: Small result sizes**

# Creating Features from Tag Ratios

**Features = Ratios?**

$$F = \left[ratio(q, t_1), \cdots, ratio(q, t_k)\right]$$

# Creating Features from Tag Ratios

Query Q

Pre-computed Tag Ratios

$Q_1$   $Q_2$   $Q_3$   $Q_4$   ...   $Q_I$

$G_1$   $G_2$   $G_3$

**Features($G_1$)   Features($G_2$)   Features($G_3$)**

With $Q_i \subseteq Q$

Group by similarity to Q

**Features based on Aggregates over ratios in a group, such as**
*SUM, AVG, STDIV, MAX, MIN*, **etc...**

# Creating Features from Tag Ratios

Canon Camera SD2
**Query Q**

Pre-computed
Tag Ratios

{Camera}
{Canon} {SD 2}
$Q_1$      $Q_2$

{Canon, Camera}
$Q_3$      $Q_4$

{Canon, Camera, SD2}
…    $Q_I$

With $Q_i \subseteq Q$

$G_1$           $G_2$           $G_3$

**Group by similarity to Q**

Features($G_1$)    Features($G_2$)   Features($G_3$)

**Features based on Aggregates over ratios in a group, such as**
*SUM, AVG, STDIV, MAX, MIN*, etc…

# Selecting queries to precompute

- $|V| > 10^7 =>$ intractable # of keyword combinations to pre-compute

**Pruning Logic**

- **Short queries:** limit query-length to $w_{max}$ words.
- **Significant correlation:**

$$ratio(q,t) \geq \Theta_{high} \frac{|\# \text{tags } t \in D|}{D} \qquad ratio(q,t) \leq \Theta_{low} \frac{|\# \text{tags } t \in D|}{D}$$

- **Ratio-support:** $|result_C(q)| \geq \alpha$



Very few keyword (combinations) satisfy *Correlation Condition*

No keyword (combination) can satisfy *Support Condition*

$f(r_w)$

$\alpha$

Number of occurrences

Rank of keyword (combination)

# Experimental Evaluation

*Task I*: **Indentifying 'Consumer-Electronics' queries**
- C = Wikipedia, T = Entity-Categories (contained in pages)
- **Accuracy:** 93.0% (n-grams only)
            93.2% (n-grams + Brand/Models/Product Type/ Product Attribute lexicons))
            95.6%  (Tag ratios only)
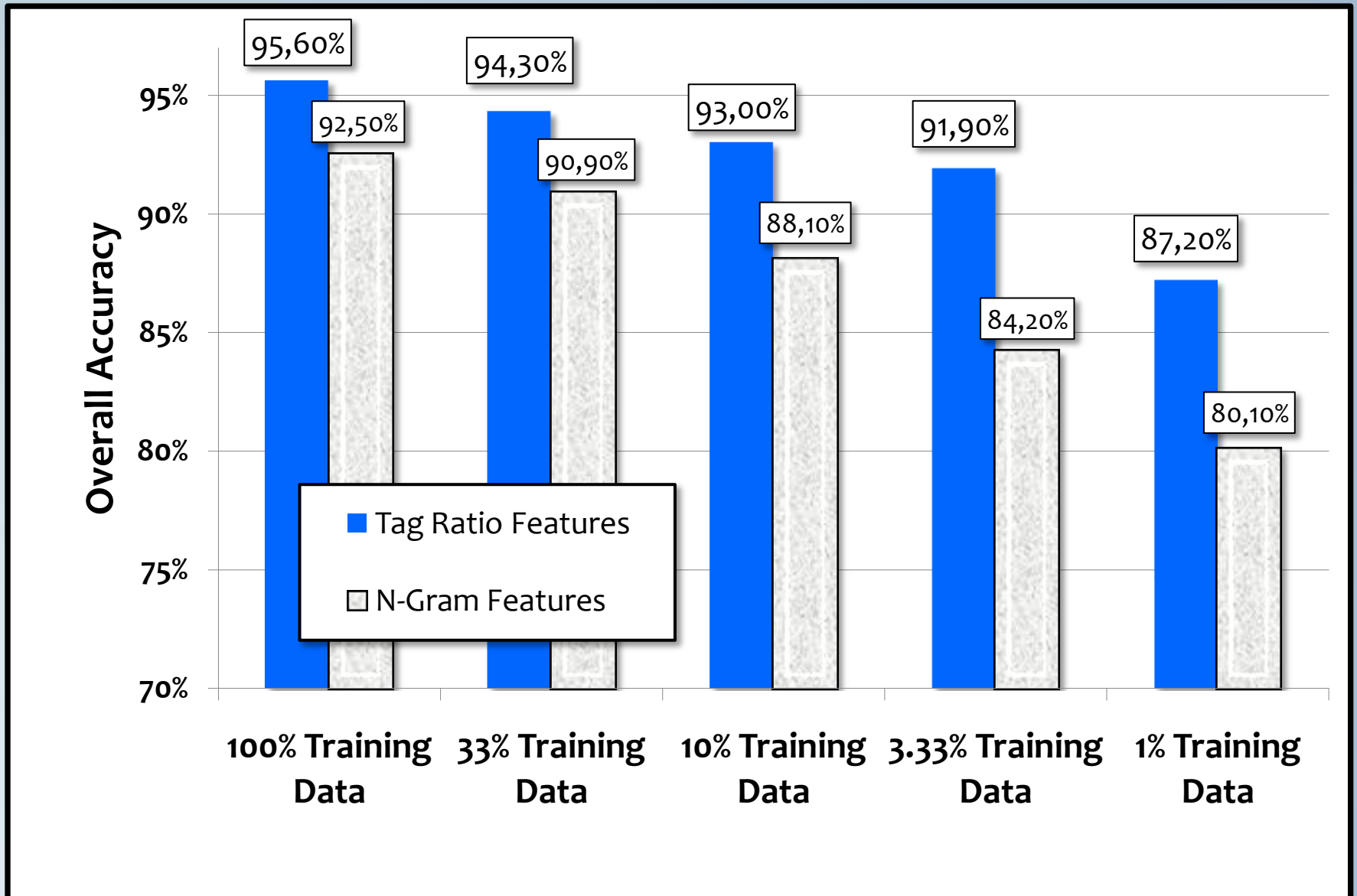            96.5%  (Tag ratios+ n-grams)

*Task II*: **Indentifying 'Retail' queries**
- C = Wikipedia, T = Top Wikipedia Categories (contained in pages)
- + C = Sponsored Search Bids, T = Advertiser IDs (top advertisers)
- Large training corpus (~330K labeled examples)
- **Accuracy:** 92.5% (n-grams only)  => 93.3%  (Tag Ratios+ ngrams)

*Task III*: **Indentifying 'Heath'-related queries**
- Same corpora/tags as before
- Very large training corpus (~800K labeled examples)
- **Accuracy:** 98.2% (n-grams only)  => 98.8%  (Tag Ratios + ngrams)

# Experimental Evaluation: Generalization

# Experimental Evaluation: Query Selection

Using earlier classification tasks, we evaluate features based on:
- Single-Word queries only
- Single-word queries + selected query/ratio combinations
- All queries in training/test data + all subsets

Results:
- Pruning results in very large reduction in space of ratios to store ($\Theta_{low}$ = 0.8, $\Theta_{high}$ = 1.2 => 0.8% of ratios (for frequent keywords) remain).
- Differences in classification accuracy slight: (0.17% or less)

# Many thanks!

# Any Questions?