## Improving Search Through *Efficient* A/B Testing:

#### **A Case Study**

Nokia Maps "Place Discovery" Team, Berlin:

Hannes Kruppa, Steffen Bickel, Mark Waldaukat, Felix Weigel, Ross Turner, Peter Siemen



### **Nokia Maps for Everyone!**

#### Nokia Maps on iPhone, iPad and Android

ZOMGitsCj.com 381 videos 😒 Subscribe





TDC Rate/24: 150 kr./md. i 24 mdr. I alt 3.600 kr. ÅOP 0 %. Kreditomkostninger 0 kr. Kontantpris u/abn. 4.399 kr.



Nokia 701 Hands-on Preview by boybandreject00 9,953 views



Nokia N8 HD Video Sample -Driving Through the ... by boybandreject00 1,669 views



Nokia 5233 Review : Camera UI by boybandreject00





NOKI

Uploader Comments (boybandreject00)

## Nokia Maps Team, Berlin





🔲 🔟 Te	elekom.de 🛜 💦 💈	<b>8</b> 🐚 22:10
Near Rue Ja	rby acob 14, Paris FRANCE	
Ţ	Ladurée 21, Rue Bonaparte 21, 75006 ★★★★	100 m Paris
<b>T</b>	<b>Tsukizi</b> Rue Des Ciseaux 2, 75006 Par	100 m is
Ţ	Le Petit Saint-Benoît Rue Saint-Benoît 4, 75006 Pa ★ ★ ★ ★	100 m ris
Ţ	Rhumerie (La) 166 Boulevard Saint-Germain	<b>100 m</b> 166, 7
Ţ	Café Mabillon 164 Boulevard Saint-Germain ★★★★	<b>200 m</b> 164, 7
Ţ	<b>Le Petula</b> 6, rue des Ciseaux 6, 75006 p	200 m aris
<b>M</b>	<b>Del Tapa</b> Rue De Buci 38, 75006 Paris	200 m
÷	Q. (b) t=	•••



























NOK

\* 22:01

### **Problem: Which Places to Show?**

- Restaurants? Hotels? Shopping? ...
- rank by Ratings?
- Distance?
- Usage?
- Trending?
- ....





### **Approach: A/B-Test Different Versions!**





## A/B-Test for *Nearby Places*

#### Version A: Best of Eat'n'Drink



#### Version B: *Best of Hotels*



Versions Compete for User engagement:

= Number of Actions performed on places.



#### There Is A Better Approach For Ranked Lists

#### [Joachims et al 2008]:

"How Does Clickthrough Data Reflect Retrieval Quality?"

- Classical A/B testing converges slowly for ranked lists
- Classical A/B testing often doesn't reflect actual relevance
- A/B Tests for Ranked Result Lists: Rank- Interleaving
- Use Rank-Interleaving for faster statistical significance



## Efficient A/B Testing: Rank Interleaving

#### Version A: Best of Eat'n'Drink

Te	elekom.de 🛜 💦 🏄 🦻	22:10		
Nearby Rue Jacob 14, Paris FRANCE				
<b>¥1</b>	Ladurée 21, Rue Bonaparte 21, 75006 Pa ★★★★	100 m ris		
<b>ĭ1</b>	<b>Tsukizi</b> Rue Des Ciseaux 2, 75006 Paris	100 m		
<b>Ť1</b>	Le Petit Saint-Benoît Rue Saint-Benoît 4, 75006 Paris ★ ★ ★ ★	100 m		
<b>ĭ1</b>	Rhumerie (La) 166 Boulevard Saint-Germain 16	100 m 6, 7		
<b>ĭ1</b>	Café Mabillon 164 Boulevard Saint-Germain 16 ★ ★ ★ ★	<b>200 m</b> 64, 7		
<b>ĭ1</b>	<b>Le Petula</b> 6, rue des Ciseaux 6, 75006 pari	200 m S		
M	<b>Del Tapa</b> Rue De Buci 38, 75006 Paris	200 m		
4	Q. 🕲 🕇 🖛	•••		

#### Version B: *Best of Hotels*





## Efficient A/B Testing: *Rank Interleaving*

#### Version A: Best of Eat'n'Drink

III T	elekom.de 🛜 👘 🐉 🔹	22:10
Neal Rue Ja	rby acob 14, Paris FRANCE	X
<b>T</b>	Ladurée 21, Rue Bonaparte 21, 75006 Pa ★ ★ ★ ★ ★	100 m Iris
<b>ĭ1</b>	<b>Tsukizi</b> Rue Des Ciseaux 2, 75006 Paris	100 m
<b>ĭ</b> 1	Le Petit Saint-Benoît Rue Saint-Benoît 4, 75006 Paris ★★★★	100 m
<b>ĭ</b> 1	<b>Rhumerie (La)</b> 166 Boulevard Saint-Germain 16	<b>100 m</b> 56, 7
<b>Y</b> 1	Café Mabillon 164 Boulevard Saint-Germain 16 ★ ★ ★ ★ ★	<b>200 m</b> 54, 7
<b>Y</b> 1	<b>Le Petula</b> 6, rue des Ciseaux 6, 75006 pari	200 m S
M	<b>Del Tapa</b> Rue De Buci 38, 75006 Paris	200 m
÷	Q. (b) t=	•••

┿

#### Version B: *Best of Hotels*



#### Rank Interleaving: *Version A* + *B*

1	Ladurée 21, Rue Bonaparte 21, 75006 Pa ★ ★ ★ ★	100 m aris
	<b>Hotel Adion</b> Unter den Linden 77, Berlin, 10	122 m 1437
	<b>Tsukizi</b> Rue Des Ciseaux 2, 75006 Paris	100 m
	<b>Hotel Merian</b> Friedrichstrasse 3, Berlin, 1043	263 m 7
1	Le Petit Saint-Benoît Rue Saint-Benoît 4, 75006 Paris * * * *	100 m
	<b>Hotel WestInn</b> Gendarmenmarkt 3, Berlin, 104	489 km 37



## **Randomized Mixing of Result Lists**

• Interleaved list is filled with pairs of results, one item from each version. Coin toss decides who comes first.



#### Interleaved Result list

<empty>



Version B 1. beta 2. kappa 3. tau













• Interleaved list is filled with pairs of results, one item from each version. Coin toss decides who comes first.



#### **Version A**

- 1. alpha
- 2. (beta)
- 3. gamma
- 4. delta
- 5. epsilon

Final list shown to user

1. alpha (from A)

- 2. beta (from B)
- 3. gamma (from A)
- 4. kappa (from B)
- 5. tau (from B)
- 6. delta (from A)
- 7. epsilon (from A, extra)



Version B 1. beta 2. kappa



### **Declaring A Winner**

- Statistical Significance Test
- Input (after hadoop-based log-processing...)
  - Number of clicks on version A
  - Number of clicks on version B
- G-Test:
  - improved version of Pearson's Chi-squared test.
  - G > 6.635 corresponds to 99% confidence level
- Null hypothesis:
  - Frequency of counts is equally distributed over both versions.
- Test statistic:

$$G = 2 \sum_{i \in \{A,B\}} [\text{counts i}] \ln \left( \frac{[\text{counts i}]}{[\text{total counts/2}]} \right)$$



## **Managing Multiple Versions**



# **Managing Multiple Versions**



#### **Caveat 1: Randomization**

- don't confuse users with changing results, i.e.: provide a consistent user experience
- Solution:
  - Random generator is seeded with USER-ID for each query.
  - Each user gets his personal random generator.



### **Caveat 2: Healthy Click Data**

- we are relying on the integrity of transmitted user actions
- sensitive to log contamination (unidentified QA, spam)
- user-clicks plot:





### Caveat 3: A/B Clicks vs. Coverage

- Coverage = non-empty responses (in percent)
- For example
  - A/B interleaving of eat&drink vs. eat&drink + going out
  - difference is not significant
  - But coverage different, percentage of responses with POIs nearby:
    - 60% eat&drink
    - 62% eat&drink + going out
- Higher coverage decides in case there is no statistical difference



#### Case Study: Eat'n'Drink versus Hotels: Not the User Behaviour we had expected!



NOKIA

#### Case Study: Eat'n'Drink versus Hotels: Not the User Behaviour we had expected!



NOKIA

#### Summary

- use A/B Rank Interleaving to optimize result relevance
- Rank Interleaving is easy to implement. Works.
- in a distributed search architecture manage your A/B test configurations conveniently using Zookeeper
- harness your hadoop/search analytics stack for A/B test evaluations
- don't make assumptions about your users!

• [Joachims et al 2008]:

"How Does Clickthrough Data Reflect Retrieval Quality?"



#### Thanks!

#### Get in touch: hannes.kruppa@nokia.com

Nokia Maps "Place Discovery" Team, Berlin:

Hannes Kruppa, Steffen Bickel, Mark Waldaukat, Felix Weigel, Ross Turner, Peter Siemen

