# Predicting Search Satisfaction Metrics with Interleaved Comparisons

Anne Schuth

Katja Hofmann

Filip Radlinski

University of Amsterdam

Microsoft
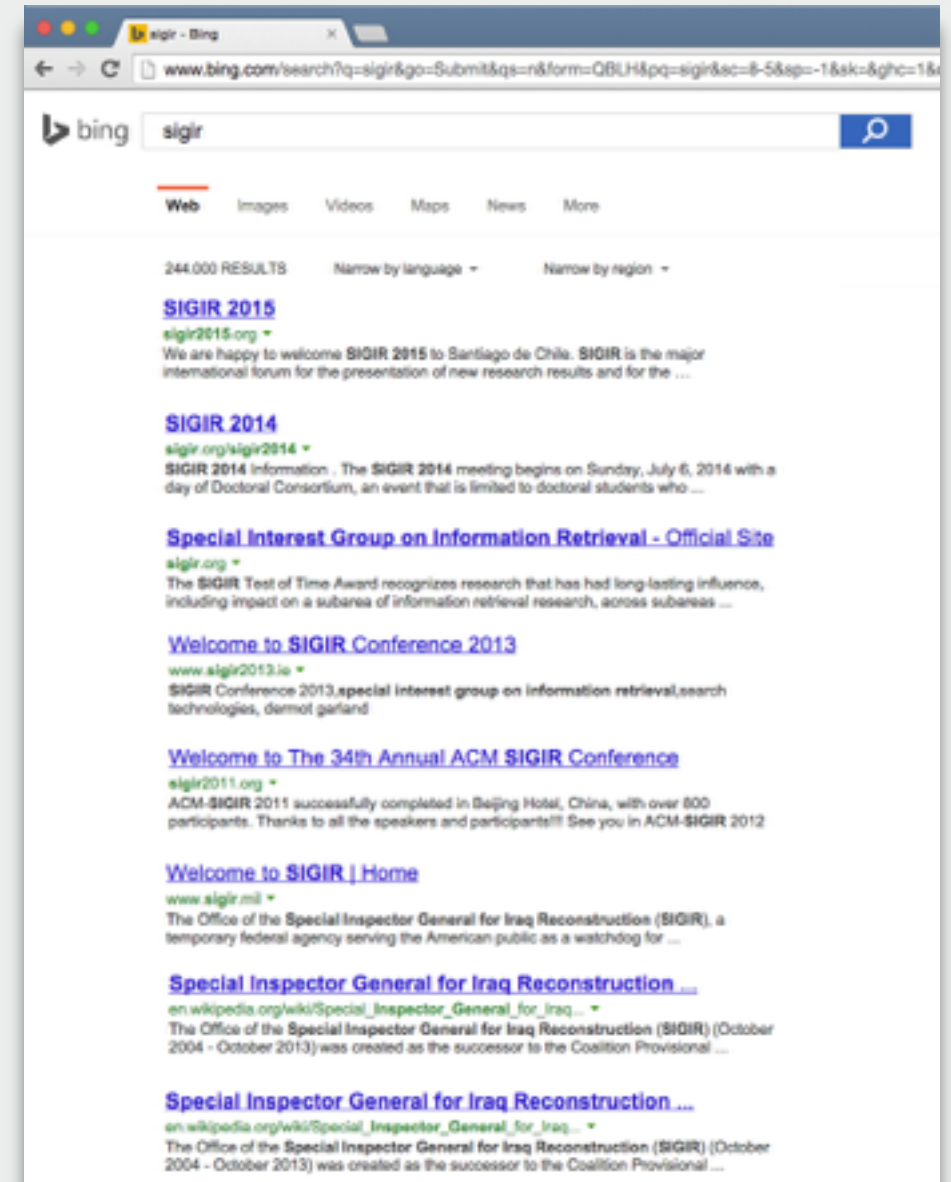
Microsoft

anne.schuth@uva.nl

katja.hofmann@microsoft.com

filiprad@microsoft.com

# Motivation - **Evaluation**



or

# Motivation - **Evaluation**
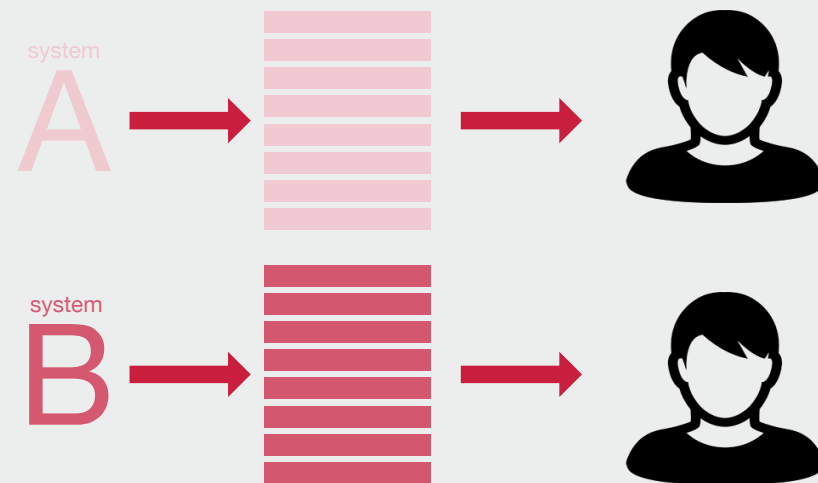
system
# A or system B

# Motivation - **AB Testing**
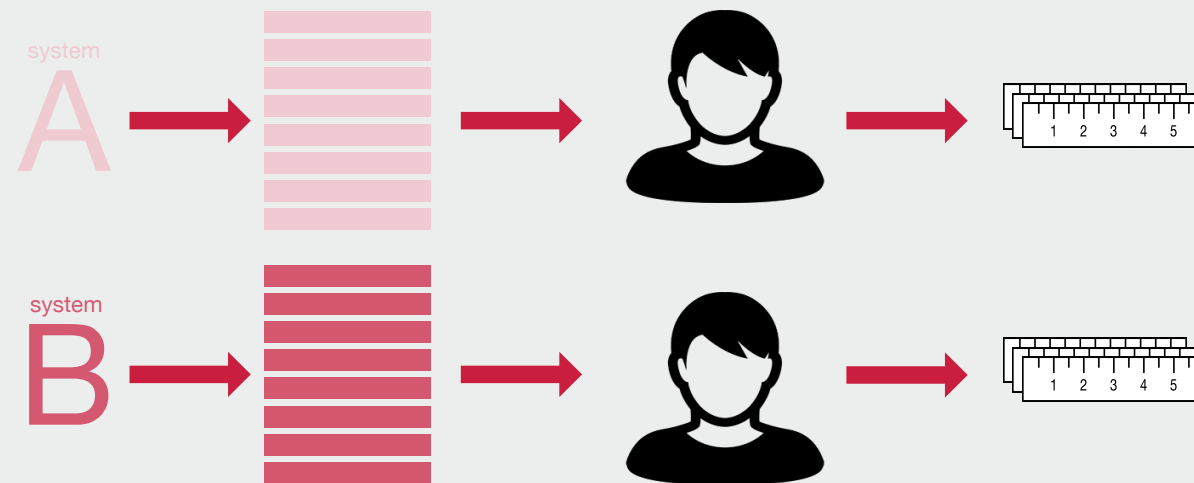
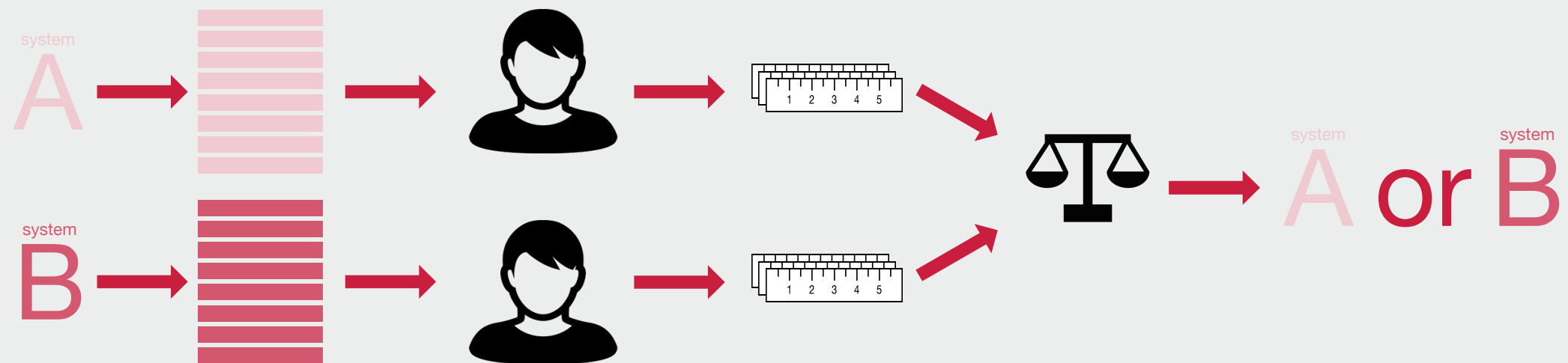system
A →

system
B →

# Motivation - **AB Testing**



✣ User population **divided** into two groups

# Motivation - **AB Testing**

❖ User population **divided** into two groups
❖ Trusted and **sophisticated metrics**

# Motivation - **AB Testing**



✤ User population **divided** into two groups

✤ Trusted and **sophisticated metrics**

✤ **Difference in metric** indicates the winner

# Motivation - **AB Testing**



✤ User population **divided** into two groups

✤ Trusted and **sophisticated metrics**

✤ **Difference in metric** indicates the winner

✤ **Between subject** design

❖ Differences between users and their queries

❖ **Low sensitivity**, millions of queries

# Motivation - **Interleaving**

system
A →

system
B →

# Motivation - **Interleaving**

✤ Users see **both** systems

# Motivation - **Interleaving**



✤ Users see **both** systems

✤ **Simple metric:** system with more clicks wins

# Motivation - **Interleaving**



- ✤ Users see **both** systems
- ✤ **Simple metric:** system with more clicks wins
- ✤ **Within subject** design
  - ❖ **Both systems** now cater for **every user**
  - ❖ **High sensitivity**, 10-100x less queries needed (compared to AB Testing)

# Motivation - **Team Draft Interleaving (TDI)**

| A | B |
|---|---|
| doc 1 | doc 2 |
| doc 2 | doc 4 |
| doc 3 | doc 7 |
| doc 4 | doc 1 |
| doc 5 | doc 3 |

F. Radlinski, M. Kurup, and T. Joachims. How does clickthrough data reflect retrieval quality? In CIKM '08. 2008

Predicting Search Satisfaction Metrics with Interleaved Comparisons

# Motivation - **Team Draft Interleaving (TDI)**

A       B

doc 1

doc 2

doc 4

doc 3

doc 7

F. Radlinski, M. Kurup, and T. Joachims. How does clickthrough data reflect retrieval quality? In CIKM '08. 2008

Predicting Search Satisfaction Metrics
with Interleaved Comparisons

5

# Motivation - **Team Draft Interleaving (TDI)**

F. Radlinski, M. Kurup, and T. Joachims. How does clickthrough data reflect retrieval quality? In CIKM '08. 2008

Predicting Search Satisfaction Metrics with Interleaved Comparisons

# Motivation - **Team Draft Interleaving (TDI)**

✤ **Infer** winner per query

 ❖ System with more **clicks** wins

 ❖ A < B

✤ Count **number of wins** over many queries

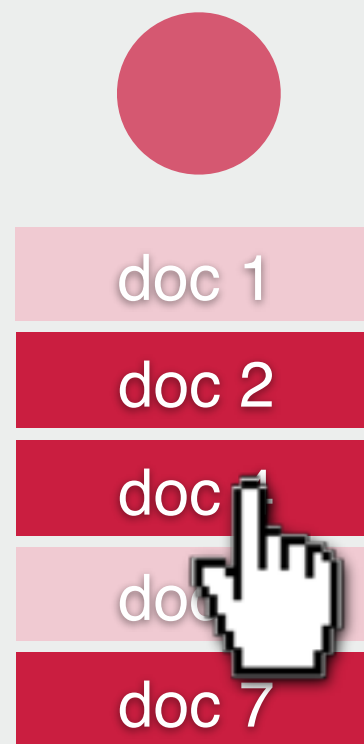| doc 1 |
| doc 2 |
| doc 1 |
| doc |
| doc 7 |

F. Radlinski, M. Kurup, and T. Joachims. How does clickthrough data reflect retrieval quality? In CIKM '08. 2008

Predicting Search Satisfaction Metrics
with Interleaved Comparisons

# Motivation - AB Testing - **As a Gold Standard**

# Motivation - AB Testing - **As a Gold Standard**



1st click, 5sec dwell time

# Motivation - AB Testing - **As a Gold Standard**



1st click, 5sec dwell time

2nd click, user stays away

# Motivation - AB Testing - **Metrics**

| AB Metric | Description |
|-----------|-------------|
|           |             |

# Motivation - AB Testing - **Metrics**

| AB Metric | Description |
|-----------|-------------|
| AB | Fraction queries with at least one **click** |

# Motivation - AB Testing - **Metrics**

| AB Metric | Description |
|-----------|-------------|
| AB | Fraction queries with at least one **click** |
| AB@1 | Fraction queries with at least one **click on 1st position** |

# Motivation - AB Testing - **Metrics**

| AB Metric | Description |
|---|---|
| AB | Fraction queries with at least one **click** |
| AB@1 | Fraction queries with at least one **click on 1st positio** |
| AB$_S$ | Fraction queries with at least one **SAT click** |

Classifier predicting **SAT probability** with a **threshold**

# Motivation - AB Testing - **Metrics**

| AB Metric | Description |
|---|---|
| AB | Fraction queries with at least one **click** |
| AB@1 | Fraction queries with at least one **click on 1st position** |
| $AB_S$ | Fraction queries with at least one **SAT click** |
| $AB_S@1$ | Fraction queries with at least one **SAT click on 1st position** |

Classifier predicting **SAT probability** with a **threshold**

# Motivation - AB Testing - **Metrics**

| AB Metric | Description |
|---|---|
| AB | Fraction queries with at least one **click** |
| AB@1 | Fraction queries with at least one **click on 1st position** |
| $AB_S$ | Fraction queries with at least one **SAT click** |
| $AB_S$@1 | Fraction queries with at least one **SAT click on 1st position** |
| $AB_T$ | **Time** from the query issue **until first click** |

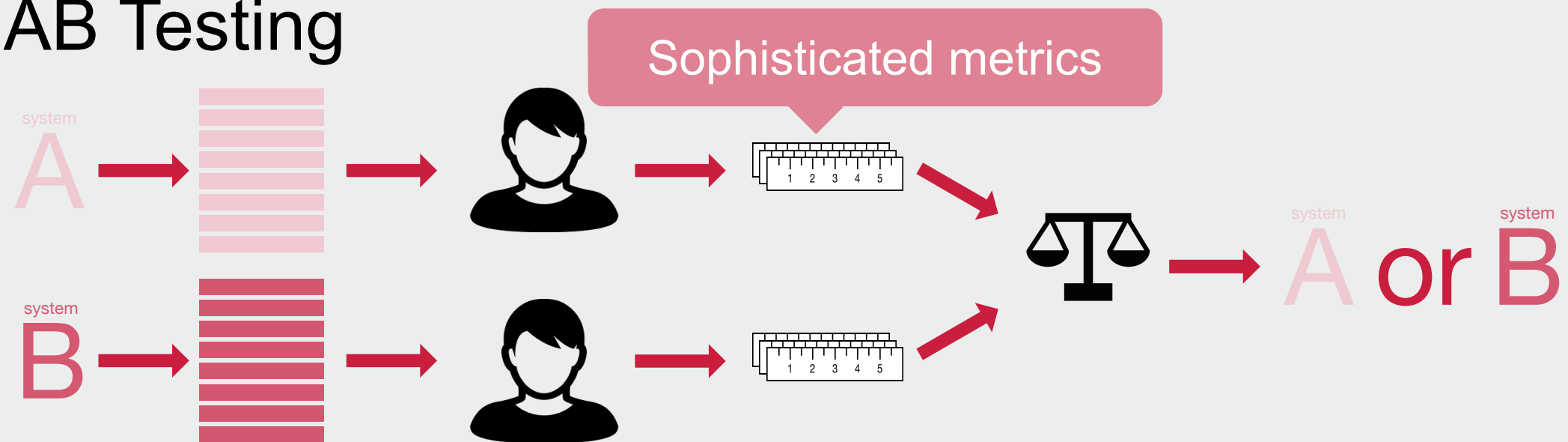Classifier predicting **SAT probability** with a **threshold**

# Motivation - AB Testing - **Metrics**

| AB Metric | Description |
|---|---|
| AB | Fraction queries with at least one **click** |
| AB@1 | Fraction queries with at least one **click on 1st position** |
| $AB_S$ | Fraction queries with at least one **SAT click** |
| $AB_S$@1 | Fraction queries with at least one **SAT click on 1st position** |
| $AB_T$ | **Time** from the query issue **until first click** |
| $AB_T$@1 | **Time** from the query issue **until first click on 1st position** |
| $AB_{T,S}$ | **Time** from the query issue **until first SAT click** |
| $AB_{T,S}$@1 | **Time** from the query issue **until first SAT click on 1st position** |

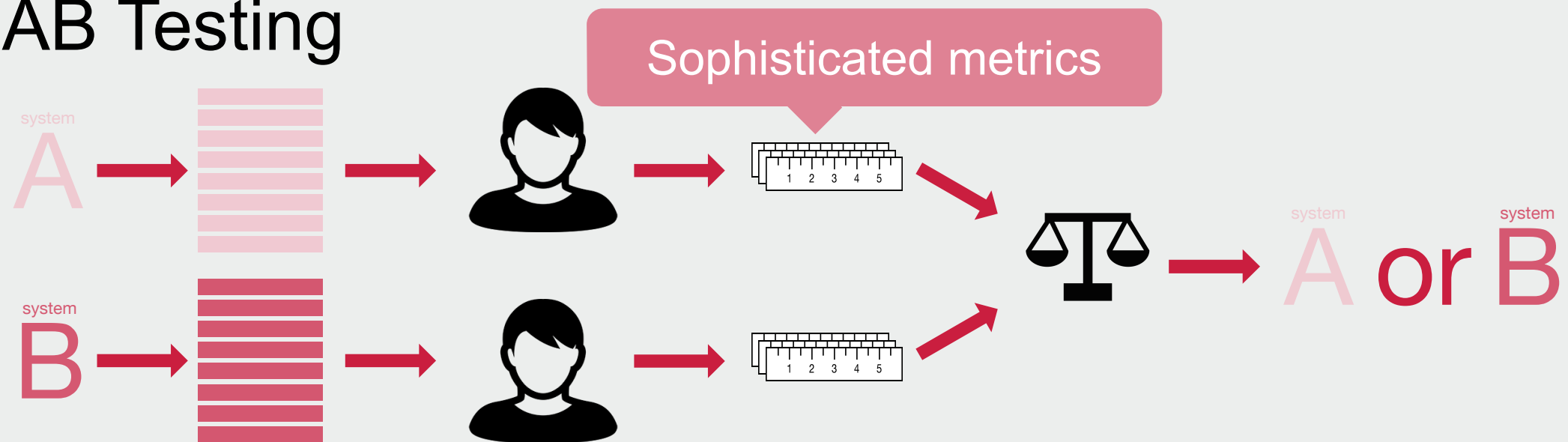Classifier predicting **SAT probability** with a **threshold**
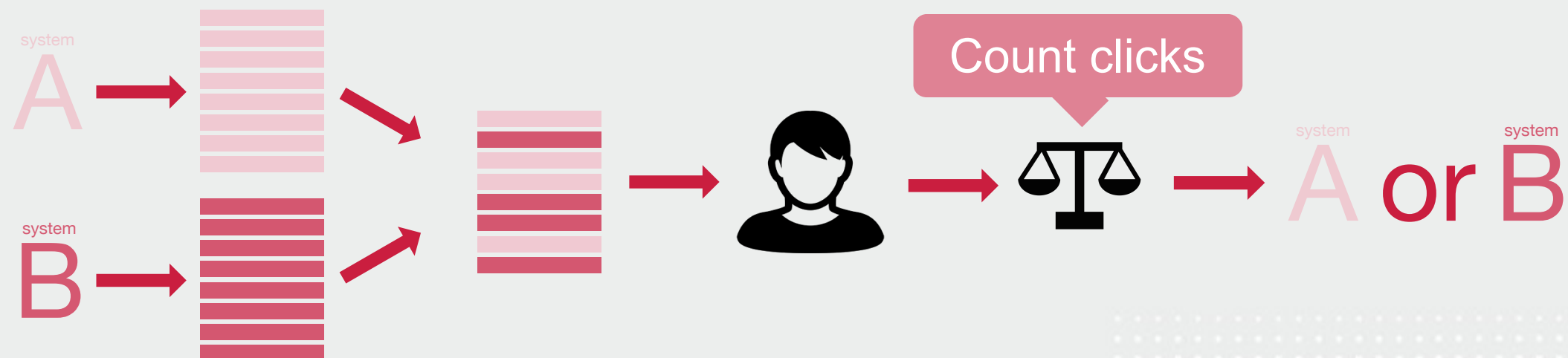
# Motivation - **Agreement**

❖ AB Testing



Sophisticated metrics

system A

system B

A or B

# Motivation - **Agreement**

✤ AB Testing

Sophisticated metrics

system A → [list] → [person] → [1 2 3 4 5] ↘

⚖ → system A or system B

system B → [list] → [person] → [1 2 3 4 5] ↗

✤ Interleaving

Count clicks

system A →
system B → [combined list] → [person] → ⚖ → system A or system B

Motivation - **Agreement**

❖ AB Testing

Sophisticated metrics

❖ Interleaving

Count clicks

# Outline

Motivation
**Data + analysis**
Methods + results
Conclusions

# Data - **Properties**

# Data - **Properties**

✤ **38 ranker pairs**

# Data - **Properties**

- ✤ **38 ranker pairs**
  - ❖ AB Tested + Interleaved (TDI)

# Data - **Properties**

✤ **38 ranker pairs**

  ❖ AB Tested + Interleaved (TDI)

  ❖ only **ranking** changes

# Data - **Properties**

- ✤ **38 ranker pairs**
  - ❖ AB Tested + Interleaved (TDI)
  - ❖ only **ranking** changes
  - ❖ bing.com, web, desktop

# Data - **Properties**

- ✤ **38 ranker pairs**
  - ❖ AB Tested + Interleaved (TDI)
  - ❖ only **ranking** changes
  - ❖ bing.com, web, desktop
  - ❖ 9 months in 2014

# Data - **Properties**

✤ **38 ranker pairs**

- ❖ AB Tested + Interleaved (TDI)
- ❖ only **ranking** changes
- ❖ bing.com, web, desktop
- ❖ 9 months in 2014
- ❖ United States locale

# Data - **Properties**

✤ **38 ranker pairs**

  ❖  AB Tested + Interleaved (TDI)

  ❖  only **ranking** changes

  ❖  bing.com, web, desktop

  ❖  9 months in 2014

  ❖  United States locale

✤ **Click volume**

# Data - **Properties**

✤ **38 ranker pairs**

   ❖ AB Tested + Interleaved (TDI)

   ❖ only **ranking** changes

   ❖ bing.com, web, desktop

   ❖ 9 months in 2014

   ❖ United States locale

✤ **Click volume**

   ❖ AB: ~1 week, **high** volume

# Data - **Properties**

- ✤ **38 ranker pairs**
  - ❖ AB Tested + Interleaved (TDI)
  - ❖ only **ranking** changes
  - ❖ bing.com, web, desktop
  - ❖ 9 months in 2014
  - ❖ United States locale
- ✤ **Click volume**
  - ❖ AB: ~1 week, **high** volume
  - ❖ Interleaving: ~4 days, **low** volume

# Data - **Properties**

- ✤ **38 ranker pairs**
  - ❖ AB Tested + Interleaved (TDI)
  - ❖ only **ranking** changes
  - ❖ bing.com, web, desktop
  - ❖ 9 months in 2014
  - ❖ United States locale
- ✤ **Click volume**
  - ❖ AB: ~1 week, **high** volume
  - ❖ Interleaving: ~4 days, **low** volume
  - ❖ **~80 times** more queries for AB

# Data - **Properties**

- ✤ **38 ranker pairs**
  - ❖ AB Tested + Interleaved (TDI)
  - ❖ only **ranking** changes
  - ❖ bing.com, web, desktop
  - ❖ 9 months in 2014
  - ❖ United States locale
- ✤ **Click volume**
  - ❖ AB: ~1 week, **high** volume
  - ❖ Interleaving: ~4 days, **low** volume
  - ❖ **~80 times** more queries for AB
  - ❖ **~3 billion clicks**

# Data - Analysis - **Agreement**

✤ **Interleaving (TDI)** does **not agree** well with **AB metrics**

| AB Metric | Interleaving (TDI) |
|---|---|
| AB | 0.63 |

# Data - Analysis - **Agreement**

✤ **Interleaving (TDI)** does **not agree** well with **AB metrics**

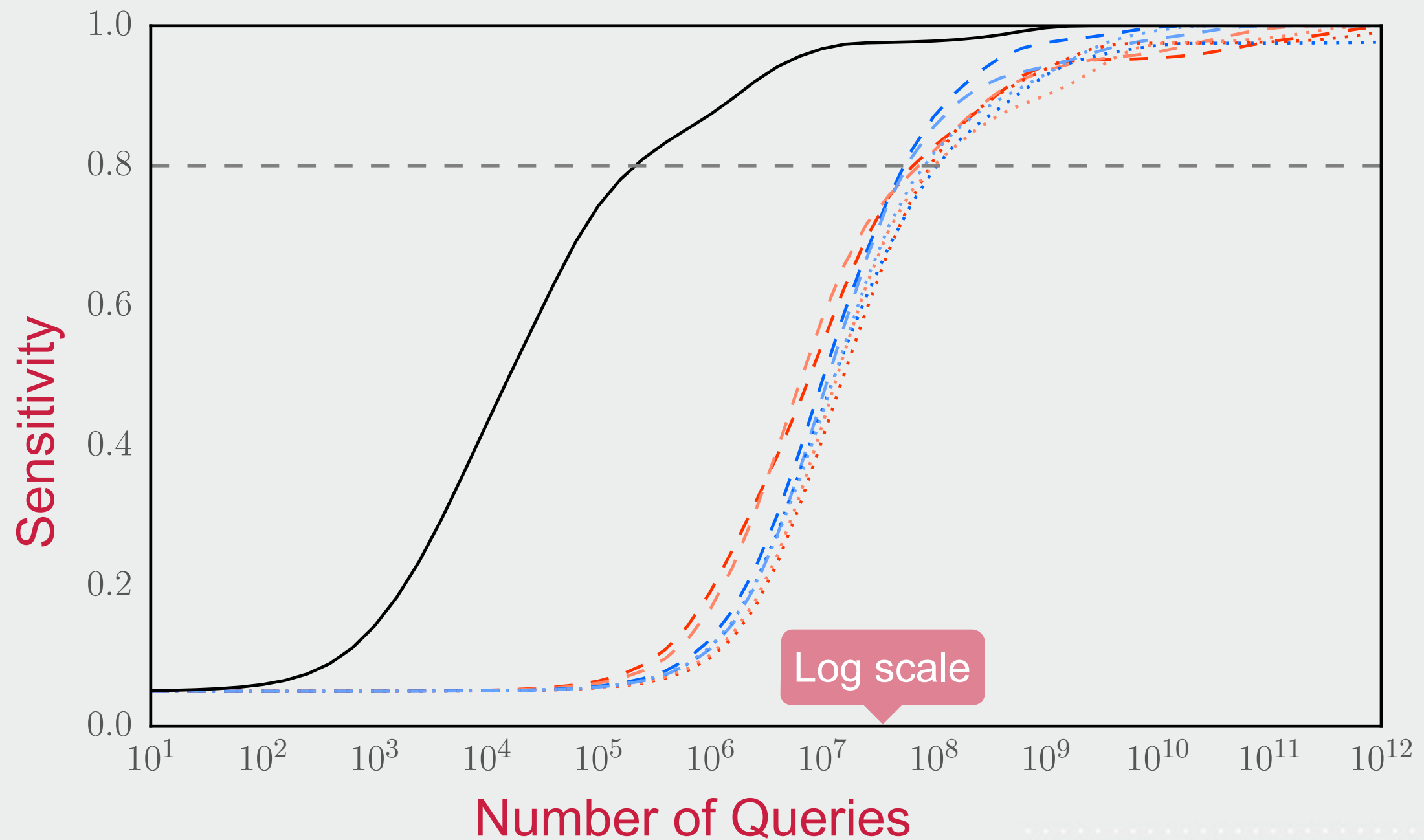| AB Metric | Interleaving (TDI) |
|---|---|
| AB | 0.63 |
| AB@1 | **0.71** |
| $AB_S$ | **0.71** |
| $AB_S$@1 | **0.76** |
| $AB_T$ | 0.53 |
| $AB_T$@1 | 0.45 |
| $AB_{T,S}$ | 0.47 |
| $AB_{T,S}$@1 | 0.42 |

Significantly different from random

# Data - Analysis - **Sensitivity (Power)**

✤ **How many queries** are required for statistically significant conclusions?

# Data - Analysis - **Sensitivity (Power)**

✤ **How many queries** are required for statistically significant conclusions?

✤ Sensitivity (power) analysis

  ❖ alpha=0.05, two sided

  ❖ AB Testing: **independent** t-test

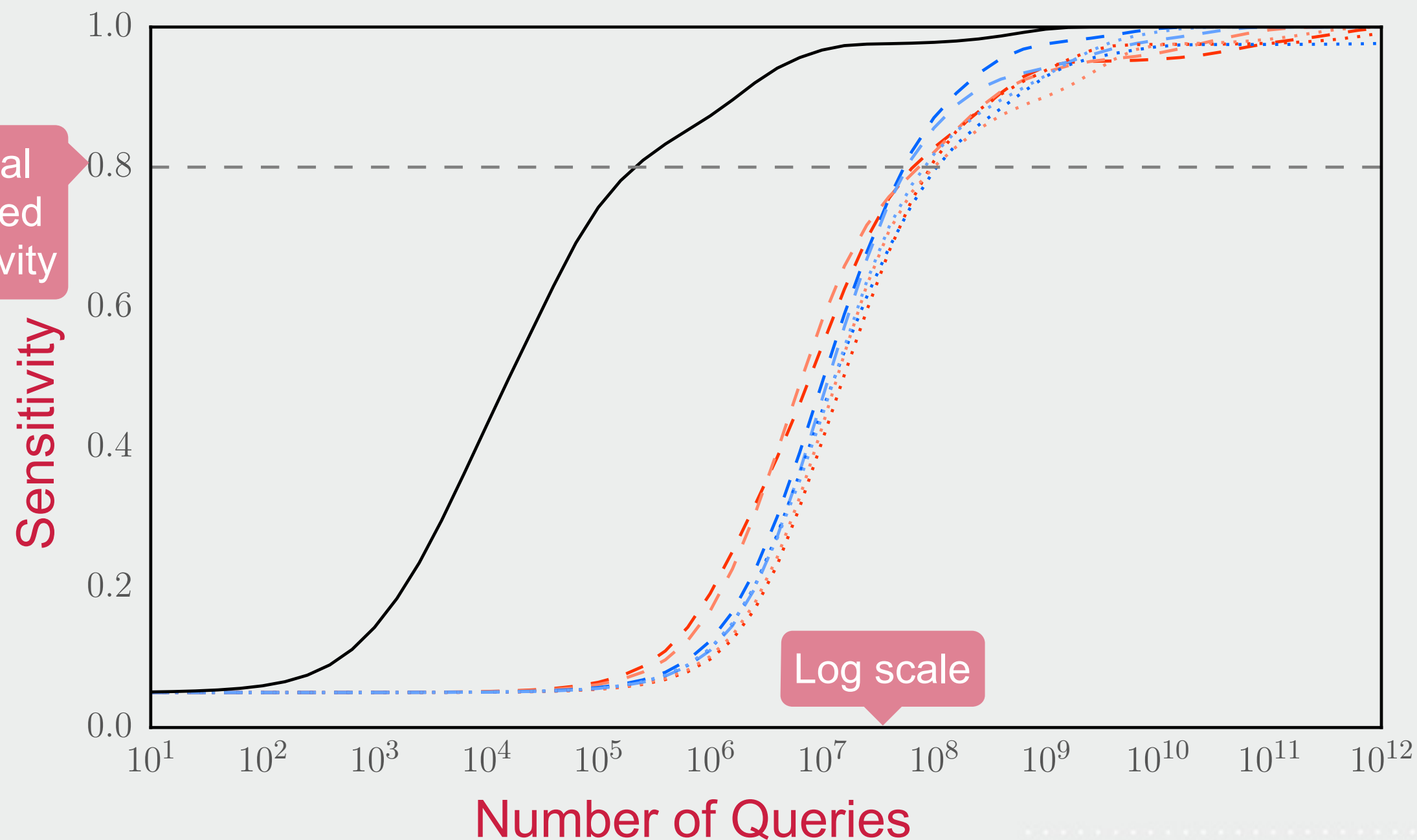  ❖ Interleaving (TDI): **paired** t-test

# Data - Analysis - **Sensitivity**
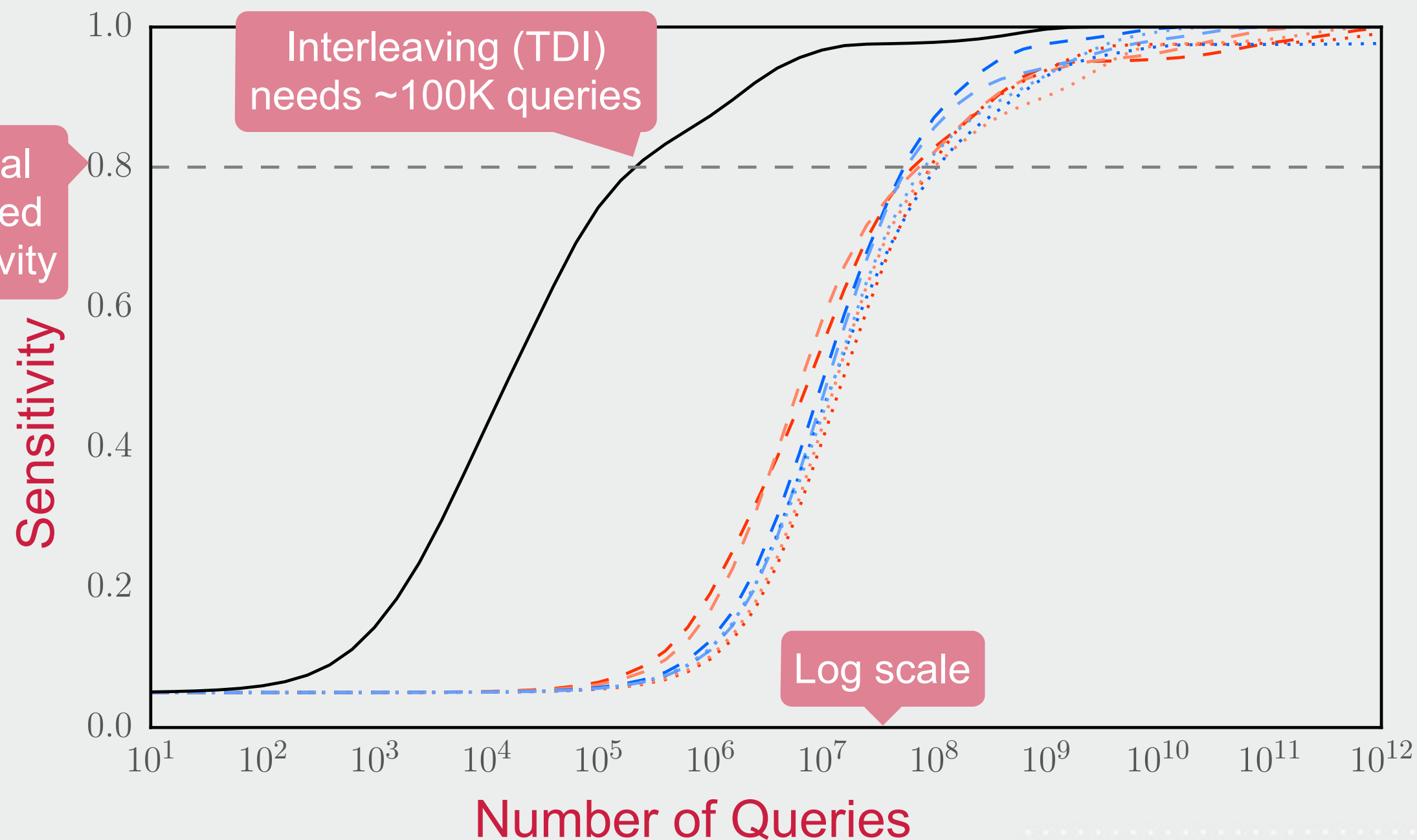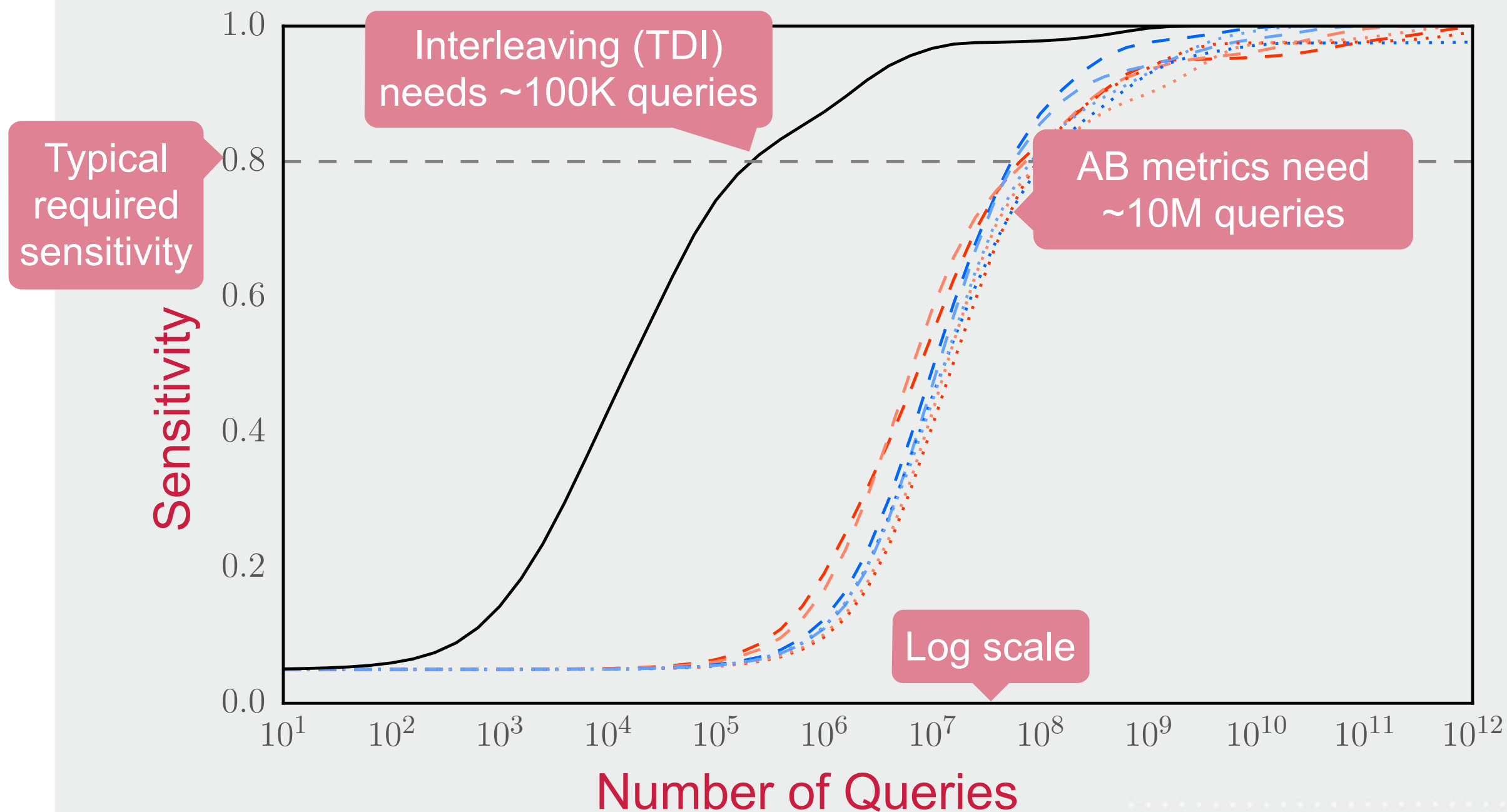
# Data - Analysis - **Sensitivity**

Log scale

# Data - Analysis - **Sensitivity**

Typical required sensitivity

Sensitivity

Log scale

Number of Queries

# Data - Analysis - **Summary**

# Data - Analysis - **Summary**

✤ **AB Testing** has **low sensitivity**

# Data - Analysis - **Summary**

✤ **AB Testing** has **low sensitivity**

✤ **Interleaving** (TDI) has **high sensitivity** (10-100x AB)

# Data - Analysis - **Summary**

✤ **AB Testing** has **low sensitivity**

✤ **Interleaving** (TDI) has **high sensitivity** (10-100x AB)

✤ **Interleaving** (TDI) has **low agreement** with AB metrics

# Data - Analysis - **Summary**

✤ **AB Testing** has **low sensitivity**

✤ **Interleaving** (TDI) has **high sensitivity** (10-100x AB)

✤ **Interleaving** (TDI) has **low agreement** with AB metrics

We aim to

Improve interleaving (TDI) to **increase agreement** with a given AB metric while **maintaining sensitivity**

# Data - Analysis - **Aim**

| | Sensitivity (required #queries) | Agreement with AB (prefer same ranker) |
|---|---|---|
| AB Testing | ~10M ✖ | ~90% ✔ |

# Data - Analysis - **Aim**

| | Sensitivity (required #queries) | Agreement with AB (prefer same ranker) |
|---|---|---|
| AB Testing | ~10M ✕ | ~90% ✓ |
| Interleaving (TDI) | ~100K ✓ | ~60% ✕ |

# Data - Analysis - **Aim**

| | Sensitivity (required #queries) | Agreement with AB (prefer same ranker) |
|---|---|---|
| **AB Testing** | ~10M ✖ | ~90% ✔ |
| **Interleaving (TDI)** | ~100K ✔ | ~60% ✖ |
| **Improved Interleaving (TDI)** | **~100K ?** ✔ | **~90% ?** ✔ |

# Outline

Motivation
Data + analysis
**Methods + results**
Conclusions

# Methods

1. **Matching AB Metrics**
2. Parameterized Credit Functions
3. Combined Credit Functions

# Methods - **Matching AB Metric**

# Methods - **Matching AB Metric**

✤ **Interleaving** traditionally counts **all clicks**

# Methods - **Matching AB Metric**

✤ **Interleaving** traditionally counts **all clicks**

✤ **Instead** of counting all clicks …

# Methods - **Matching AB Metric**

✤ **Interleaving** traditionally counts **all clicks**

✤ **Instead** of counting all clicks …

✤ … we propose to **match AB metrics**

# Methods - **Matching AB Metric**

✤ **Interleaving** traditionally counts **all clicks**

✤ **Instead** of counting all clicks …

✤ … we propose to **match AB metrics**

   ❖ Count only **certain** clicks

# Methods - **Matching AB Metric**

✤ **Interleaving** traditionally counts **all clicks**

✤ **Instead** of counting all clicks …

✤ … we propose to **match AB metrics**

> ❖ Count only **certain** clicks
>> ✦ @1

# Methods - **Matching AB Metric**

✤ **Interleaving** traditionally counts **all clicks**

✤ **Instead** of counting all clicks …

✤ … we propose to **match AB metrics**

    ❖   Count only **certain** clicks

        ✦  @1

        ✦  SAT

# Methods - Matching AB Metric - **Sensitivity**

Methods - Matching AB Metric - **Sensitivity**

# Methods - Matching AB Metric - **Sensitivity**

Methods - Matching AB Metric - **Sensitivity**

# Methods - Matching AB metric - **Agreement**

Vanilla interleaving

| @ | TDI | TDI@1 | $TDI_S$ | $TDI_S@1$ | $TDI_T$ | $TDI_T@1$ | $TDI_{T,S}$ | $TDI_{T,S}@1$ |
|---|---|---|---|---|---|---|---|---|
| **AB** | 0.63 | | | | | | | |
| **AB@1** | 0.71 | 0.68 | | | | | | |
| **$AB_S$** | 0.71 | | **0.87** | | | | | |
| **$AB_S@1$** | 0.76 | | | 0.63 | | | | |
| **$AB_T$** | 0.53 | | | | **0.71** | | | |
| **$AB_T@1$** | 0.45 | | | | | 0.58 | | |
| **$AB_{T,S}$** | 0.47 | | | | | | 0.58 | |
| **$AB_{T,S}@1$** | 0.42 | | | | | | | 0.58 |

# Methods - Matching AB metric - **Agreement**

Vanilla interleaving

| | TDI | TDI@1 | $TDI_S$ | $TDI_S$@1 | $TDI_T$ | $TDI_T$@1 | $TDI_{T,S}$ | $TDI_{T,S}$@1 |
|---|---|---|---|---|---|---|---|---|
| **AB** | 0.63 | 0.66 | **0.84** | 0.66 | 0.61 | 0.61 | 0.58 | 0.53 |
| **AB@1** | 0.71 | 0.68 | **0.76** | 0.63 | 0.63 | 0.47 | 0.55 | 0.55 |
| **$AB_S$** | 0.71 | 0.68 | **0.87** | 0.68 | 0.68 | 0.58 | 0.61 | 0.55 |
| **$AB_S$@1** | 0.76 | 0.68 | **0.82** | 0.63 | 0.74 | 0.53 | 0.61 | 0.50 |
| **$AB_T$** | 0.53 | 0.55 | 0.47 | 0.55 | **0.71** | 0.55 | 0.68 | 0.58 |
| **$AB_T$@1** | 0.45 | 0.47 | 0.45 | 0.58 | **0.63** | 0.58 | 0.61 | 0.62 |
| **$AB_{T,S}$** | 0.47 | 0.55 | 0.53 | **0.71** | 0.66 | 0.66 | 0.58 | 0.53 |
| **$AB_{T,S}$@1** | 0.42 | 0.50 | 0.53 | **0.66** | 0.61 | **0.66** | 0.58 | 0.58 |

# Methods - Matching AB metric - **Agreement**

Vanilla interleaving

| | TDI | TDI@1 | $TDI_S$ | $TDI_S$@1 | $TDI_T$ | $TDI_T$@1 | $TDI_{T,S}$ | $TDI_{T,S}$@1 |
|---|---|---|---|---|---|---|---|---|
| **AB** | 0.63 | 0.66 | **0.84** | 0.66 | 0.61 | 0.61 | 0.58 | 0.53 |
| **AB@1** | 0.71 | 0.68 | **0.76** | 0.63 | 0.63 | 0.47 | 0.55 | 0.55 |
| **$AB_S$** | 0.71 | 0.68 | **0.87** | 0.68 | 0.68 | 0.58 | 0.61 | 0.55 |
| **$AB_S$@1** | 0.76 | 0.68 | **0.82** | 0.63 | 0.74 | 0.53 | 0.61 | 0.50 |
| **$AB_T$** | 0.53 | 0.55 | 0.47 | 0.55 | **0.71** | 0.55 | 0.68 | 0.58 |
| **$AB_T$@1** | 0.45 | 0.47 | 0.45 | 0.58 | **0.63** | 0.58 | 0.61 | 0.62 |
| **$AB_{T,S}$** | 0.47 | 0.55 | 0.53 | **0.71** | 0.66 | 0.66 | 0.58 | 0.53 |
| **$AB_{T,S}$@1** | 0.42 | 0.50 | 0.53 | **0.66** | 0.61 | **0.66** | 0.58 | 0.58 |

Highest agreement not on diagonal

Predicting Search Satisfaction Metrics
with Interleaved Comparisons

20

# Methods

1. Matching AB Metrics
2. **Parameterized Credit Functions**
3. Combined Credit Functions

# Methods - **Parametrized Credit**

# Methods - **Parametrized Credit**

✤ We aim to increase agreement

# Methods - **Parametrized Credit**

✤ We aim to increase agreement

Remember, we have a model that predicts **SAT probability**

✤ **Parameterize TDI** with a SAT threshold $t_s$

❖ $TDI_S{}^{ts}$ and $TDI_{T,S}{}^{ts}$

# Methods - **Parametrized Credit**

Remember, we have a model that predicts **SAT probability**

✤ We aim to increase agreement

✤ **Parameterize TDI** with a SAT threshold $t_s$

❖ $TDI_S^{ts}$ and $TDI_{T,S}^{ts}$

Click based    Time based

# Methods - **Parametrized Credit**

❖ We aim to increase agreement

Remember, we have a model that predicts **SAT probability**

❖ **Parameterize TDI** with a SAT threshold $t_s$

❖ $TDI_S^{ts}$ and $TDI_{T,S}^{ts}$

Click based

Time based

Filter out non SAT clicks, **can reduce** sensitivity

# Methods - **Parametrized Credit**

❖ We aim to increase agreement

*Remember, we have a model that predicts **SAT probability***

❖ **Parameterize TDI** with a SAT threshold $t_s$

❖ $TDI_S{}^{ts}$ and $TDI_{T,S}{}^{ts}$

*Click based*  *Time based*

*Filter out non SAT clicks, **can reduce** sensitivity*

❖ Find **optimal threshold $t_s$**

❖ Maximize agreement for **each** AB metric

❖ Repeat n=100 times:

❖ Take bootstrap sample

❖ Grid search to find $t_s$ that maximizes agreement

❖ Report performance on "out of bag" sample

# Methods - Parametrized Credit - **Sensitivity**

# Methods - Parametrized Credit - **Sensitivity**

# Methods - Parametrized Credit - **Agreement**

| AB Metric | Vanilla | Click based |
|---|---|---|
| | TDI | TDI$_S^{ts}$ |
| AB | 0.63 | **0.82** |
| AB@1 | 0.71 | |
| AB$_S$ | 0.71 | |
| AB$_S$@1 | 0.76 | |
| AB$_T$ | 0.53 | |
| AB$_T$@1 | 0.45 | |
| AB$_{T,S}$ | 0.47 | |
| AB$_{T,S}$@1 | 0.42 | |

# Methods - Parametrized Credit - **Agreement**

| AB Metric | Vanilla TDI | Click based $\text{TDI}_\text{s}^\text{ts}$ |
|---|---|---|
| AB | 0.63 | **0.82** |
| AB@1 | 0.71 | **0.79** |
| $\text{AB}_\text{S}$ | 0.71 | **0.84** |
| $\text{AB}_\text{S}$@1 | 0.76 | **0.84** |
| $\text{AB}_\text{T}$ | 0.53 | 0.47 |
| $\text{AB}_\text{T}$@1 | 0.45 | **0.49** |
| $\text{AB}_\text{T,s}$ | 0.47 | 0.46 |
| $\text{AB}_\text{T,s}$@1 | 0.42 | 0.52 |

# Methods - Parametrized Credit - **Agreement**

| AB Metric | Vanilla TDI | Click based $TDI_S^{ts}$ | Time based $TDI_{T,S}^{ts}$ |
|---|---|---|---|
| AB | 0.63 | **0.82** | 0.53 |
| AB@1 | 0.71 | **0.79** | 0.54 |
| $AB_S$ | 0.71 | **0.84** | 0.48 |
| $AB_S$@1 | 0.76 | **0.84** | 0.48 |
| $AB_T$ | 0.53 | 0.47 | **0.67** |
| $AB_T$@1 | 0.45 | **0.49** | **0.62** |
| $AB_{T,S}$ | 0.47 | 0.46 | **0.61** |
| $AB_{T,S}$@1 | 0.42 | 0.52 | **0.62** |

# Methods - Parametrized Credit - **Agreement**

| AB Metric | Vanilla | Click based | Time based |
|---|---|---|---|
| | TDI | $TDI_S^{ts}$ | $TDI_{T,S}^{ts}$ |
| AB | 0.63 | **0.82** | 0.53 |
| AB@1 | 0.71 | **0.79** | 0.54 |
| $AB_S$ | 0.71 | **0.84** | 0.48 |
| $AB_S$@1 | 0.76 | **0.84** | 0.48 |
| $AB_T$ | 0.53 | 0.47 | **0.67** |
| $AB_T$@1 | 0.45 | **0.49** | **0.62** |
| $AB_{T,S}$ | 0.47 | 0.46 | **0.61** |
| $AB_{T,S}$@1 | 0.42 | 0.52 | **0.62** |

# Methods

1. Matching AB Metrics
2. Parameterized Credit Functions
3. **Combined Credit Functions**

# Methods - **Combined Credit**

# Methods - **Combined Credit**

✤ **Combine** parameterized **credit functions**

❖    $w_S \cdot TDI_S^{ts} + w_T \cdot TDI_{T,S}^{ts}$

Click weight

Time weight

# Methods - **Combined Credit**

✤ **Combine** parameterized **credit functions**

❖ $w_S \cdot TDI_S^{ts} + w_T \cdot TDI_{T,S}^{ts}$

Click weight

Time weight

✤ Find optimal weights

❖ Maximizing agreement

# Methods - Combined Credit - **Agreement**

| AB Metric | TDI |
|---|---|
| AB | 0.63 |
| AB@1 | 0.71 |
| $AB_S$ | 0.71 |
| $AB_S$@1 | 0.76 |
| $AB_T$ | 0.53 |
| $AB_T$@1 | 0.45 |
| $AB_{T,S}$ | 0.47 |
| $AB_{T,S}$@1 | 0.42 |

# Methods - Combined Credit - **Agreement**

| AB Metric | TDI | $\text{TDI}_{T,s}^W$ agreement | Click weight $w_S$ | Time weight $w_T$ |
|-----------|-----|------------|------|------|
| AB | 0.63 | **0.84** | *1.00* | *0.00* |
| AB@1 | 0.71 | | | |
| $\text{AB}_S$ | 0.71 | | | |
| $\text{AB}_S$@1 | 0.76 | | | |
| $\text{AB}_T$ | 0.53 | | | |
| $\text{AB}_T$@1 | 0.45 | | | |
| $\text{AB}_{T,s}$ | 0.47 | | | |
| $\text{AB}_{T,s}$@1 | 0.42 | | | |

Predicting Search Satisfaction Metrics
with Interleaved Comparisons

# Methods - Combined Credit - **Agreement**

| AB Metric | TDI | TDI$_{T,S}^W$ agreement | Click weight $w_S$ | Time weight $w_T$ |
|---|---|---|---|---|
| AB | 0.63 | **0.84** | 1.00 | 0.00 |
| AB@1 | 0.71 | **0.75** | 1.00 | 0.05 |
| AB$_S$ | 0.71 | **0.85** | 1.00 | 0.00 |
| AB$_S$@1 | 0.76 | **0.83** | 1.00 | 0.02 |
| AB$_T$ | 0.53 | **0.68** | 0.99 | 0.90 |
| AB$_T$@1 | 0.45 | **0.56** | 0.96 | 0.79 |
| AB$_{T,S}$ | 0.47 | **0.63** | 0.91 | 0.88 |
| AB$_{T,S}$@1 | 0.42 | **0.50** | 0.06 | 0.25 |

# Methods - Combined Credit - **Agreement**

| AB Metric | TDI | $\text{TDI}_{T,s}^{W}$ agreement | Click weight $w_S$ | Time weight $w_T$ |
|---|---|---|---|---|
| AB | 0.63 | **0.84** | 1.00 | 0.00 |
| AB@1 | 0.71 | **0.75** | 1.00 | 0.05 |
| AB$_S$ | 0.71 | **0.85** | 1.00 | 0.00 |
| AB$_S$@1 | 0.76 | **0.83** | 1.00 | 0.02 |
| AB$_T$ | 0.53 | **0.68** | 0.99 | 0.90 |
| AB$_T$@1 | 0.45 | **0.56** | 0.96 | 0.79 |
| AB$_{T,s}$ | 0.47 | **0.63** | 0.91 | 0.88 |
| AB$_{T,s}$@1 | 0.42 | **0.50** | 0.06 | 0.25 |

# Methods - Combined Credit - **Agreement**

| AB Metric | TDI | $\text{TDI}_{T,S}^W$ agreement | Click weight $w_S$ | Time weight $w_T$ |
|---|---|---|---|---|
| AB | 0.63 | **0.84** | 1.00 | 0.00 |
| AB@1 | 0.71 | **0.75** | 1.00 | 0.05 |
| AB$_S$ | 0.71 | **0.85** | 1.00 | 0.00 |
| AB$_S$@1 | 0.76 | **0.83** | 1.00 | 0.02 |
| AB$_T$ | 0.53 | **0.68** | 0.99 | 0.90 |
| AB$_T$@1 | 0.45 | **0.56** | 0.96 | 0.79 |
| AB$_{T,S}$ | 0.47 | **0.63** | 0.91 | 0.88 |
| AB$_{T,S}$@1 | 0.42 | **0.50** | 0.06 | 0.25 |

# Methods - Combined Credit - **Agreement**

| AB Metric | TDI | $TDI_{T,s}^W$ agreement | Click weight $w_s$ | Time weight $w_T$ |
|---|---|---|---|---|
| AB | 0.63 | **0.84** | 1.00 | 0.00 |
| AB@1 | 0.71 | **0.75** | 1.00 | 0.05 |
| $AB_s$ | 0.71 | **0.85** | 1.00 | 0.00 |
| $AB_s$@1 | 0.76 | **0.83** | 1.00 | 0.02 |
| $AB_T$ | 0.53 | **0.68** | 0.99 | 0.90 |
| $AB_T$@1 | 0.45 | **0.56** | 0.96 | 0.79 |
| $AB_{T,s}$ | 0.47 | **0.63** | 0.91 | 0.88 |
| $AB_{T,s}$@1 | 0.42 | **0.50** | 0.06 | 0.25 |

All significantly better

# Methods - Combined Credit - **Sensitivity**

Methods - Combined Credit - **Sensitivity**

Predicting Search Satisfaction Metrics with Interleaved Comparisons

28

# Methods - Combined Credit - **Sensitivity**

# Outline

Motivation
Data + analysis
Methods + results
**Conclusions**

# Conclusions - **Data Analysis**

# Conclusions - **Data Analysis**

✤ Sensitivity:

> Confirming earlier findings

❖ **AB Testing** is 10-100x **less sensitive than Interleaving**

Predicting Search Satisfaction Metrics
with Interleaved Comparisons

30

# Conclusions - **Data Analysis**

✤ Sensitivity:

Confirming earlier findings

❖ **AB Testing** is 10-100x **less sensitive than Interleaving**

✤ Agreement

New insight

❖ **Between AB Testing** and **Interleaving** (TDI) is **low:** <76%

# Conclusions - **Methods**

# Conclusions - **Methods**

✤ Interleaving (TDI) with credit **matching** AB metrics
  ❖ **Unpredictable**

# Conclusions - **Methods**

✤ Interleaving (TDI) with credit **matching** AB metrics

   ❖ **Unpredictable**

✤ Interleaving (TDI) with **parameterized** credit functions

   ❖ Improvements for **some** AB metrics

# Conclusions - **Methods**

✤ Interleaving (TDI) with credit **matching** AB metrics

    ❖ **Unpredictable**

✤ Interleaving (TDI) with **parameterized** credit functions

    ❖ Improvements for **some** AB metrics

✤ Interleaving (TDI) with **combined** credit functions

    ❖ Improvements for **all** AB metrics

# Conclusions - **Future Work**

# Conclusions - **Future Work**

✤ Consider **even richer user signals** (sessions, task level features)

# Conclusions - **Future Work**

✤ Consider **even richer user signals** (sessions, task level features)

✤ Take **magnitude** and **uncertainty** of AB metric differences into account

# Conclusions - **Future Work**

✤ Consider **even richer user signals** (sessions, task level features)

✤ Take **magnitude** and **uncertainty** of AB metric differences into account

✤ Understanding of **where and why agreement is low or high**

# Take Away

# Take Away

✤ **Rich** user **signals** in **interleaving**

# Take Away

- ✤ **Rich** user **signals** in **interleaving**
- ✤ **Agreement** of Interleaving with an AB metric can be made as high as **87%**

# Take Away

- ✤ **Rich** user **signals** in **interleaving**
- ✤ **Agreement** of Interleaving with an AB metric can be made as high as **87%**
- ✤ While maintaining **high sensitivity** of Interleaving

# Take Away

✤ **Rich** user **signals** in **interleaving**

✤ **Agreement** of Interleaving with an AB metric can be made as high as **87%**

✤ While maintaining **high sensitivity** of Interleaving

✤ Microsoft® Research

✤ UNIVERSITY OF AMSTERDAM

✤ http://anneschuth.nl

✤ @anneschuth