# Intern Presentation

A/B Testing by Interleaving

Sida Wang

# My Project

- Evaluating search relevance by interleaving results and collecting user data
  - Interleaving Framework
    - Generic, Extensible
  - Experiments to evaluate relevance by interleaving
    - Based on the paper [How Does Clickthrough Data Reflect Retrieval Quality?](#) by F. Radlinski et al

# Evaluating Search Relevance

- Without Interleaving
  - Full time human judges -> precision, recall, NDCG
  - Compare Search

# Compare Search



Result S1

Result S2

Result S3

Result G1

Result G2

Result G3

| Results from O12 trained model | Results from O14 trained model |
| --- | --- |



Sandbox  Search Results:                                REDMOND/d-siwang ▾ ❼

All Sites | People

xbox360 chip                                    🔍 Advanced

Refine Results:              Results 1-10 of about 414.                    People Matches:

Any Result Type           **Xbox 360 - Kit Oficial Brasil**          Julia Purtell
Adobe PDF                 Descritivo **Xbox360.doc** ... Imagens: fornecidas pelo Gerente      THERMAL
Email                     de Conta Microsoft ... e capacidade para até quatro controles;     TEST
Excel                     **chip** gráfico ATI 500 MHz; três processadores simétricos .. Sinta   ENGINEER
                          a ação do jogo com a ..
show more ▾               Authors: hsaraiva  Date: 11/8/2007  Size: 103KB
                          http://sharepoint/sites/brz_x360/Games Description/Descritivo           John Van
Any Site                  **Xbox360**.doc - Explain Rank                            Ness
                                                                                    VLSI
portals/vote              **Xbox360** Silicon Operations:                           PRODUCT
sharepoint/.../audi...    with Supply Chain partners Ease of integration **Xbox360** Silicon   ENGINEER
my/.../kevkelly           Ops Solution SAP Adaptor: .. not let them just happen. Enable
sharepoint/.../vote       Agility Loosely couple processes and solutions. **Xbox360** ..        Norm
                          Authors: Robert Meshew  Date: 3/14/2006  Size: 3MB        LeMieux
show more ▾               http://team/sites/SRM/Shared                            HARDWARE
                          Documents/SiliconOps_SMART.ppt  View duplicates - Explain   DESIGN
                          Rank                                                     ENGINEER/
Any Author                                                                         TEST
                          **Xbox360** Case Study
Kenneth Bell (Volt)       New product New channels New supplier model New financial   View More People »
v-pumuru                  considerations **Xbox360** - The Vision .. Inventory Mgmt Faster
Michael B Grossman        cycle time **Xbox360** Architecture **Xbox360** - The Reality
Rob Harris                **Xbox360** Case ..
                          Authors: Charles Fitzgerald  Date: 3/23/2006  Size: 1MB
show more ▾               http://team/sites/SRM/Shared Documents/Xbox case study2.ppt -
                          Explain Rank
Any Modified Date
                          Getting Started with XbOx 360
Since Yesterday           Just to be clear **xBoX360** is not only for games. It's part of our
Past Month                digital lifestyle and you can .. □ **xBoX360** is a Media Center
Past Six Months           aware and extender. Meaning you can connect you your console
Past Year
                          Authors: waelk  Date: 9/25/2007  Size: 1MB
Over a year ago           http://myemea/sites/waelk/Shared Documents/Getting Started
                          with XbOx 360.pdf - Explain Rank           feedback

Sandbox  Search Results:                                REDMOND/d-siwang ▾ ❼

All Sites | People

xbox360 chip                                    🔍 Advanced

Refine Results:              Results 1-10 of about 414.                    People Matches:

Any Result Type           **Xbox360** Silicon Operations:           Julia Purtell
Adobe PDF                 with Supply Chain partners Ease of integration **Xbox360** Silicon   THERMAL
Email                     Ops Solution SAP Adaptor: .. not let them just happen. Enable   TEST
Excel                     Agility Loosely couple processes and solutions. **Xbox360** ..   ENGINEER
                          Authors: Robert Meshew  Date: 3/14/2006  Size: 3MB
show more ▾               http://team/sites/SRM/Shared
                          Documents/SiliconOps_SMART.ppt  View duplicates - Explain   John Van
Any Site                  Rank                                                     Ness
                                                                                    VLSI
portals/vote              **Xbox360** Case Study                                    PRODUCT
sharepoint/.../audi...    New product New channels New supplier model New financial   ENGINEER
my/.../kevkelly           considerations **Xbox360** - The Vision .. Inventory Mgmt Faster
sharepoint/.../vote       cycle time **Xbox360** Architecture **Xbox360** - The Reality   Norm
                          **Xbox360** Case ..                                        LeMieux
show more ▾               Authors: Charles Fitzgerald  Date: 3/23/2006  Size: 1MB   HARDWARE
                          http://team/sites/SRM/Shared Documents/Xbox case study2.ppt -   DESIGN
Any Author                Explain Rank                                             ENGINEER/
                                                                                    TEST
Kenneth Bell (Volt)       **Xbox 360 - Kit Oficial Brasil**
v-pumuru                  Descritivo **Xbox360.doc** ... Imagens: fornecidas pelo Gerente   View More People »
Rob Harris                de Conta Microsoft ... e capacidade para até quatro controles;
Michael B Grossman        **chip** gráfico ATI 500 MHz; três processadores simétricos .. Sinta
                          a ação do jogo com a ..
show more ▾               Authors: hsaraiva  Date: 11/8/2007  Size: 103KB
                          http://sharepoint/sites/brz_x360/Games Description/Descritivo
Any Modified Date         **Xbox360**.doc - Explain Rank

Since Yesterday           Getting Started with XbOx 360
Past Month                Just to be clear **xBoX360** is not only for games. It's part of our
Past Six Months           digital lifestyle and you can .. □ **xBoX360** is a Media Center
Past Year                 aware and extender. Meaning you can connect you your console

Over a year ago           Authors: waelk  Date: 9/25/2007  Size: 1MB
                          http://myemea/sites/waelk/Shared Documents/Getting Started
                          with XbOx 360.pdf - Explain Rank           feedback
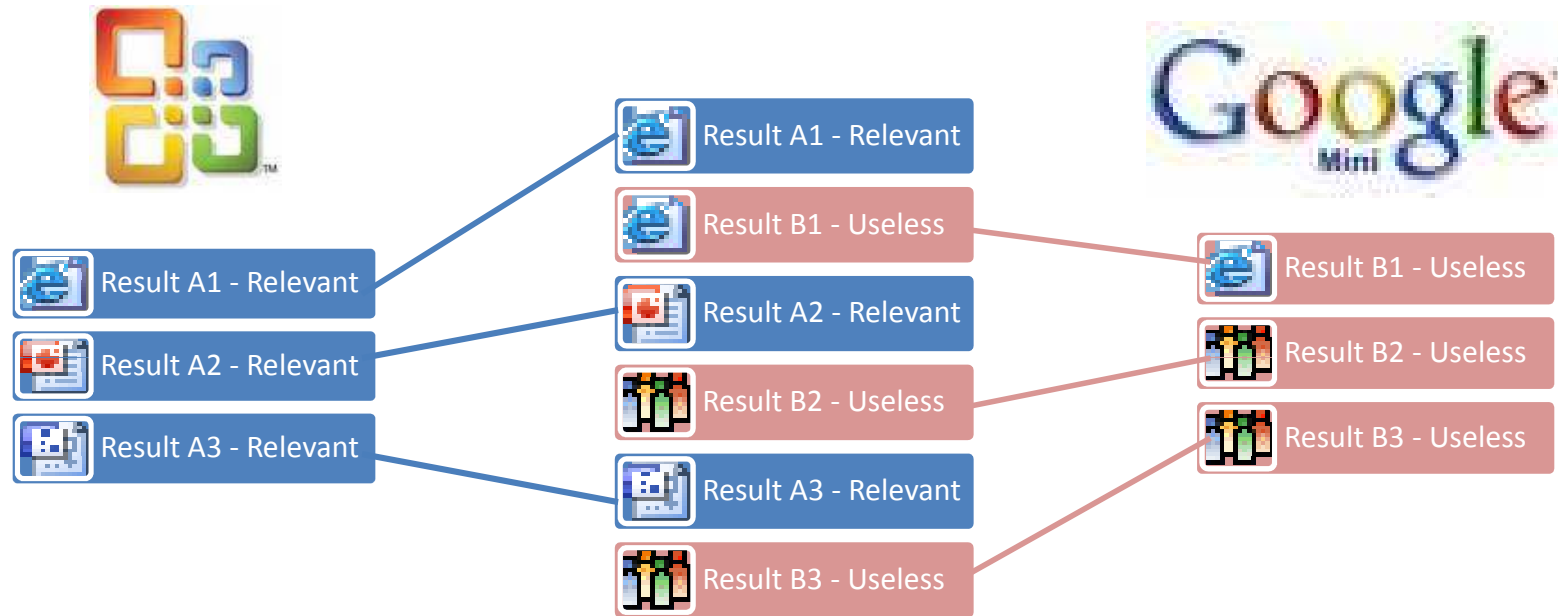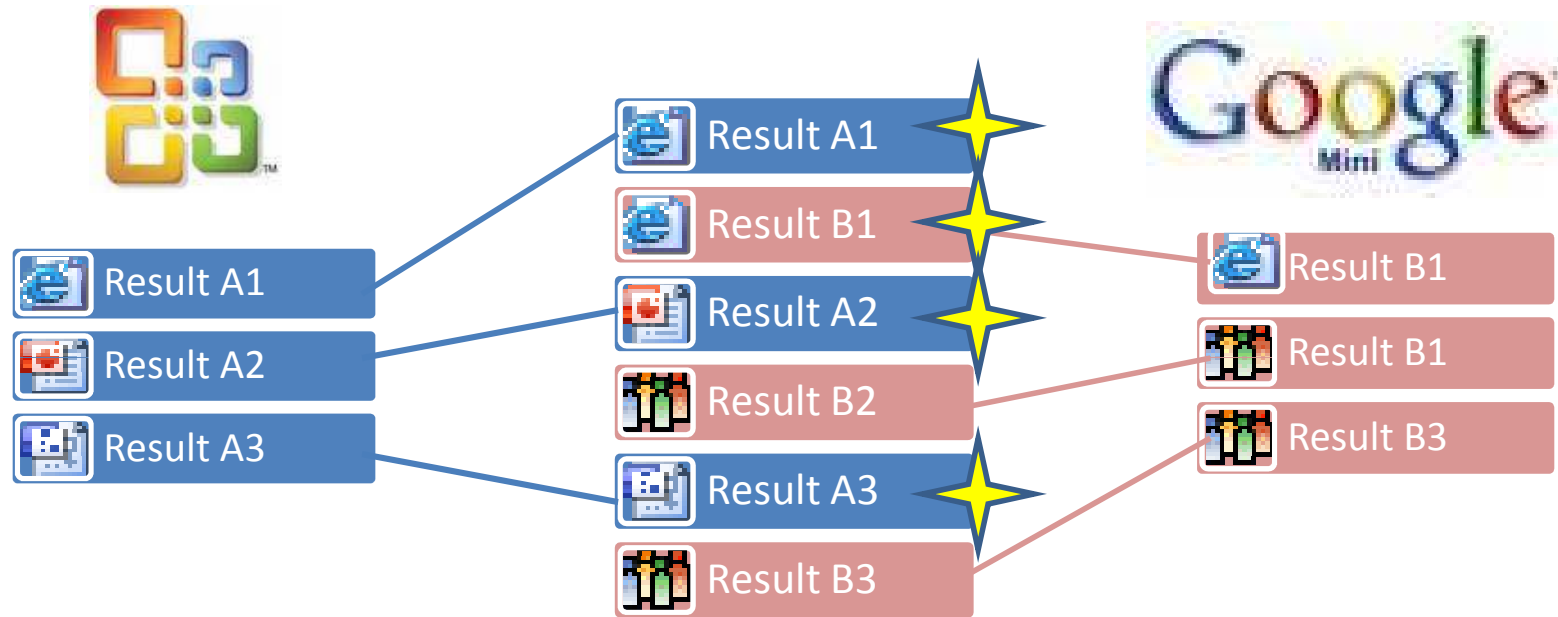
# Issues



- But do Microsoft people pick O14 Search or Google Mini?
- Maybe people tend to pick the left?
- **Alters the search experience**
  – Can never collect a lot of data using this method

# By Interleaving

# By Interleaving
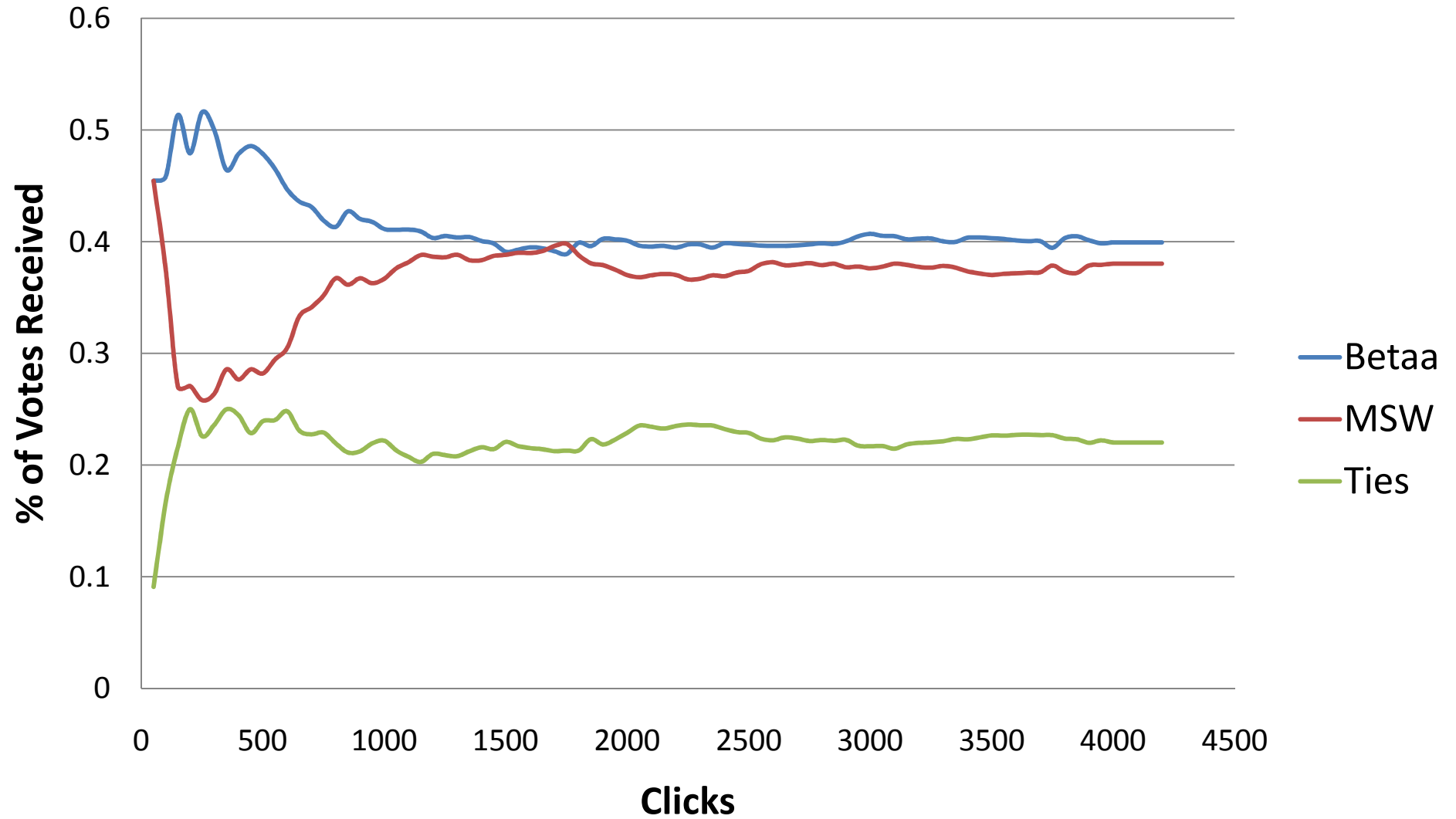
# Considerations

- Minimize impact to UX
  - So no demo, it looks exactly like normal search
- Minimize Bias
  - Summary normalization
  - Interleaving algorithms
- Reliability / performance / and the usual

# Experiments I did

- **Automated random clicks**
- Automated clicks according to relevance judgments
- Clicks from real people

# Random Clicks
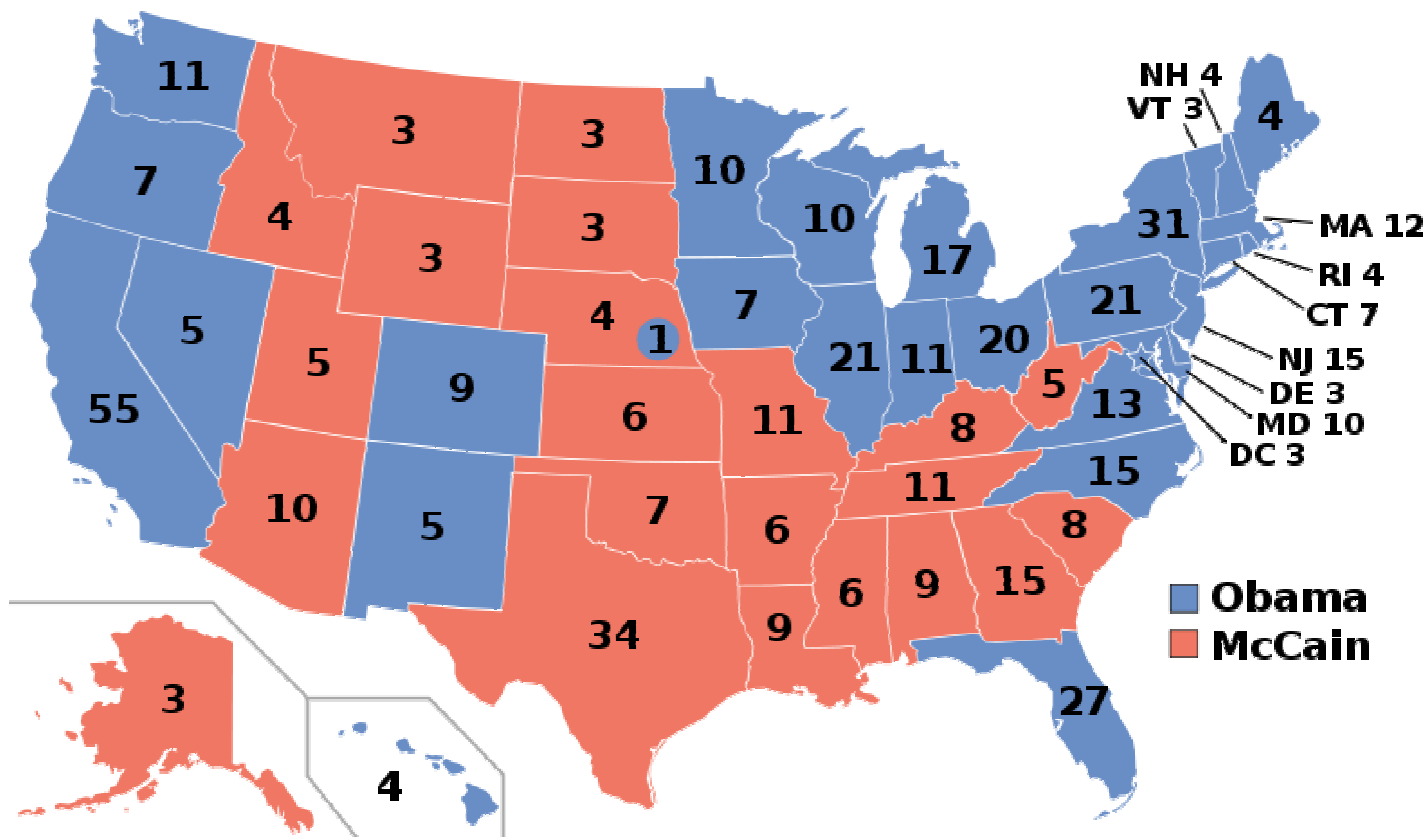


Control Using Automated Random Clicks

# A Lot of Random Clicks

## Control Using Automated Random Clicks

# Experiments I did

- Automated random clicks
- **Automated clicks according to relevance judgments**
- Clicks from real people

# O12 vs. O14

## Automated Clicks Using Relevence Judgments

- Acing05 Degraded
- Acing05
- Ties

# Experiments I did

- Automated random clicks
- Automated clicks according to relevance judgments
- **Clicks from real people**

# O12 vs. O14

## O12 vs. O14 in BSG ALL



**% of Votes Received** (y-axis)

**Clicks** (x-axis)

Legend:
- O12
- O14
- Tie

# Method of Analysis (election)

- Vote by query, by user, by session etc.
- query = person, user = state

# Summary of Results

**Method of Voting**                                    **O12 vs. O14**

by queries (direct election):        12 vs. 24

by users (1 vote per state):         4 vs. 9

by sessions (~electoral votes):      5 vs. 11

- System does not seem to matter much, but too little clicks (85) to draw significant conclusion

# What Logically Follows

- Google Mini vs. O14 (after fixing Google Mini)
- FAST vs. O14 (after fixing RSS in fssearchoffice)
- I'd love to see the results

# What can interleaving do?

- **Give relevance team more confidence**
- Use interleaving for displaying results
- Use interleaving to automatically tune the search engine

Ambition

# Add Confidence

- In addition to very traditional measures like NDCG, Precision and Recall. It is nice to have another independent metric.

- Automatic
  - Does not require human judgments

- Scalable
  - Small impact to UX

# What can interleaving do?

- Give relevance team more confidence
- **Use interleaving for displaying results**
- Use interleaving to automatically tune the search engine

Ambition

# Display

# Display

**Jimi Hendrix** – Discover music, videos, concerts, & pictures at Last.fm
Watch videos & listen to **Jimi Hendrix**: All Along the Watchtower, Purple Haze & more, plus
210 pictures. James Marshall **Hendrix** (November 27,1942 – September ...
www.last.fm/music/**Jimi+Hendrix** - Cached - Similar - 💬 🔝 ✕

Image results for **jimi hendrix** - Report images

Video results for **jimi hendrix**

**Jimi Hendrix** - All Along The
Watchtower Live ...
4 min 31 sec
www.youtube.com

**Jimi Hendrix** Purple Haze
2 min 30 sec
www.youtube.com

**Jimi** Henrix- Live at Woodstock
'69
57 min
video.google.com

**Jimi Hendrix**: The Uncut Story
- Episode 1 ...
52 min
video.google.com

News results for **jimi hendrix**

Erotic Madonna tapes, **Hendrix** contract in NY sale - 2 days ago
By Christine Kearney NEW YORK (Reuters) - Rock legend **Jimi Hendrix's** first recording
contract worth $1 and erotic audio and video tapes sent by Madonna to ...
Reuters - 496 related articles »
**Jimi Hendrix** hailed as Flying V legend - Gear4music.com - 2 related articles »

A **Jimi Hendrix** Experience
A **Jimi Hendrix** Experience, with pictures, lyrics, and music from the greatest guitar player to
ever grace the earth.
www.musicfanclubs.org/**jimihendrix**/ - Cached - Similar - 💬 🔝 ✕

# What can we do?

- Give relevance team more confidence
- Use interleave for displaying results
- **Use interleaving to automatically tune the search engine**

Ambition

# Automatic Tuning

- Many relevance models, each is good for a particular type of corpora (specs, user data, academic articles, product catalog, websites)

- Use interleaving in 10% of searches

- Use user click data to:
  - Automatically and dynamically decide on the best model, or tweak model parameters

# Thank you!

- Dmitriy, Eugene, Puneet

- Jamie, Jessica, Ping, Victor, Relevance Team

- Russ, Jon

- Search Team

- Hope to see you again in the future!

# Extra Slides

# Automatic Tuning – Pair wise?

- Pair wise comparisons scales poorly
- But there seems to be "strong stochastic transitivity"
  - Given locations A, B ,C
  - If A > B > C then $\Delta_{AC} > Max(\Delta_{AB}, \Delta_{BC})$

# How to Interleave

- Balanced
- Team Draft

# Balanced Interleaving

**Algorithm 1** Balanced Interleaving

**Input**: Rankings $A = (a_1, a_2, \ldots)$ and $B = (b_1, b_2, \ldots)$
$I \leftarrow ();\, k_a \leftarrow 1;\, k_b \leftarrow 1;$
$AFirst \leftarrow RandBit()$ .... *decide which ranking gets priority*
**while** $(k_a \leq |A|) \wedge (k_b \leq |B|)$ **do** ... *if not at end of A or B*
  **if** $(k_a < k_b) \vee ((k_a = k_b) \wedge (AFirst = 1))$ **then**
    **if** $A[k_a] \notin I$ **then** $I \leftarrow I + A[k_a]$ .. *append next A result*
    $k_a \leftarrow k_a + 1$
  **else**
    **if** $B[k_b] \notin I$ **then** $I \leftarrow I + B[k_b]$ .. *append next B result*
    $k_b \leftarrow k_b + 1$
  **end if**
**end while**
**Output**: Interleaved ranking $I$

More formally, denote $A = (a_1, a_2, \ldots)$, $B = (b_1, b_2, \ldots)$, $I = (i_1, i_2, \ldots)$, and let $c_1, c_2, \ldots$ be the ranks of the clicks w.r.t. $I$. To estimate $l$, [13] proposes to use the lowest ranked click, namely $l \approx c_{max} = \max\{c_1, c_2, \ldots\}$. Furthermore, to derive a preference between $A$ and $B$, one compares the number of clicks in the top

$$k = \min\{j : (i_{c_{max}} = a_j) \vee (i_{c_{max}} = b_j)\} \tag{1}$$

results of $A$ and $B$. In particular, the number $h_a$ of clicks attributed to $A$ and the number $h_b$ of clicks attributed to $B$ is computed as

$$h_a = |\{c_j : i_{c_j} \in (a_1, \ldots, a_k)\}| \tag{2}$$
$$h_b = |\{c_j : i_{c_j} \in (b_1, \ldots, b_k)\}|. \tag{3}$$

If $h_a > h_b$ we infer a preference for $A$, if $h_a < h_b$ we infer a preference for $B$, and if $h_a = h_b$ we infer no preference.

# A/B Testing By Interleaving



Result A1 - Relevant
Result A2 - Relevant
Result A3 - Relevant

Result A1 - Relevant
Result B1 - Useless
Result B2 - Useless
Result A3 - Relevant
Result B3 - Useless
A2 - Relevant

Result B1 - Useless
Result B2 - Useless
Result B3 - Useless