
Rebuttal Testimony of Professor Douglas W. Oard

U.S. et al. v. Google LLC

United States District Court for the District of Columbia

November 15, 2023

REDACTED FOR PUBLIC FILING

Ex. No.

UPXD105

1:20-cv-03010-APM

My Assignment

To provide my expert opinion of the analysis and opinions offered by Google's expert, Prof. Edward A. Fox, in his June 3, 2022 expert report (the "Fox Report").

REDACTED FOR PUBLIC FILING

My Overall Conclusion

Prof. Fox substantially understates the beneficial effects of user-side data on search quality.

REDACTED FOR PUBLIC FILING

Prof. Fox's Assignment

Prof. Fox states he was asked by Google counsel to:

“test the extent to which Google’s search quality is affected by the volume of user interaction data available to train its ranking algorithms”

REDACTED FOR PUBLIC FILING

Prof. Fox's Conclusions

Vast majority of Google-Microsoft search quality gap must be explained by factors other than volume of user interaction data

A company as efficient as Google could have search quality similar to Google even at Microsoft's scale

A company as efficient as Google but with Microsoft's scale would not meaningfully benefit from increase in user interaction data

There are diminishing returns to search quality from an increase in the quantity of user interaction data

REDACTED FOR PUBLIC FILING

Prof. Fox's Conclusions

Vast majority of Google-Microsoft search quality gap must be explained by factors other than volume of user interaction data

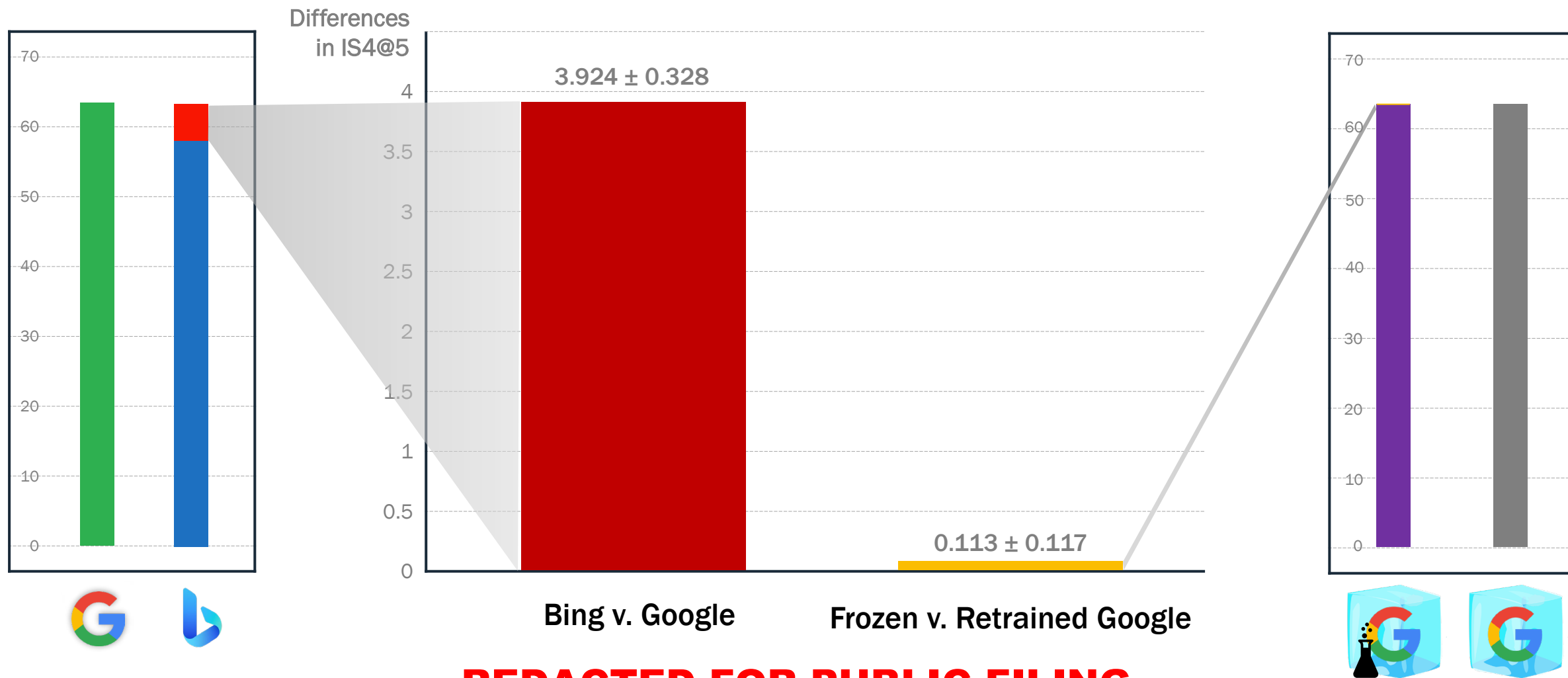
A company as efficient as Google could have search quality similar to Google even at Microsoft's scale

A company as efficient as Google but with Microsoft's scale would not meaningfully benefit from increase in user interaction data

There are diminishing returns to search quality from an increase in the quantity of user interaction data

REDACTED FOR PUBLIC FILING

The Basis for Prof. Fox's Central Conclusion



My Response to Prof. Fox's Central Conclusions

Prof. Fox's conclusions are unsupported because of:

- Unmeasured benefits of user-side data in this experiment;
- Measurement errors in the “quality gap”; and
- Important benefits of user-side data that this experiment cannot measure

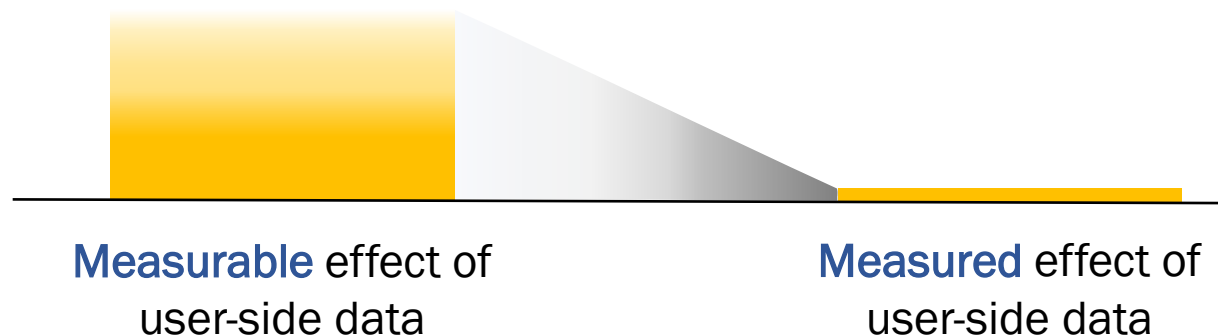
REDACTED FOR PUBLIC FILING

Many Components Not Retrained

Google **only retrained 6 components**, chosen based on their expected effect on web ranking (i.e., 10 blue links)

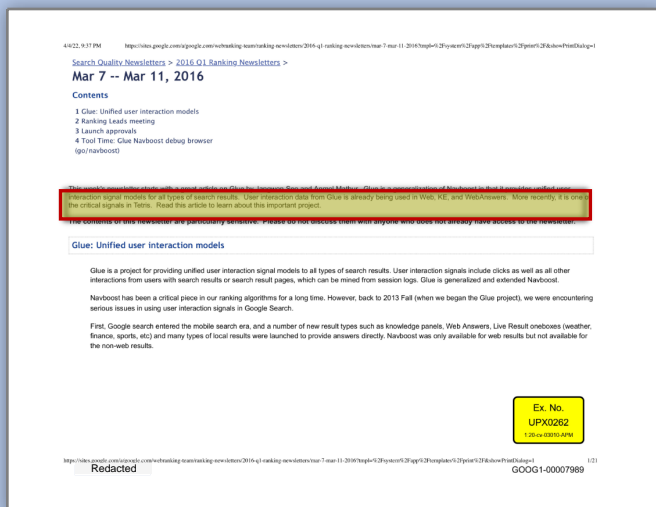
Components were not chosen based on their effect on:

- Indexing
- Spelling Correction
- Search features like images, video...
- Search advertising
- Whole-page ranking



REDACTED FOR PUBLIC FILING

Glue's Importance to Whole-Page Ranking



“User interaction data from Glue is already being used in Web, KE, and WebAnswers. More recently, it is one of the critical signals in Tetris.”

2016

REDACTED FOR PUBLIC FILING

Glue Is Used to “Trigger” and Position Search Features



Prof. Edward
Fox
Google’s Expert
Witness

“In simpler terms, **Glue aggregates diverse types of user interactions—such as clicks, hovers, scrolls, and swipes**—and creates a common metric to compare web results and search features. This process determines **both whether a search feature is triggered and where it triggers on the page.**”

REDACTED FOR PUBLIC FILING

Prof. Fox Has Never Stated that Glue Was Retrained

The Six Ranking Components Retrained in the DRE

Component Name		Component Name	
Navboost		Navboost	
RankBrain		RankBrain	
DeepRank		DeepRank	
QBST		QBST	
Term Weighting		Term Weighting	
RankEmbed-BERT		RankEmbedBERT	

Expert Report of Edward A. Fox, App. A at Table 1

Google
DXD-26.004

REDACTED FOR PUBLIC FILING

The IS4@5 Metric Evaluates Web and Search Features Results

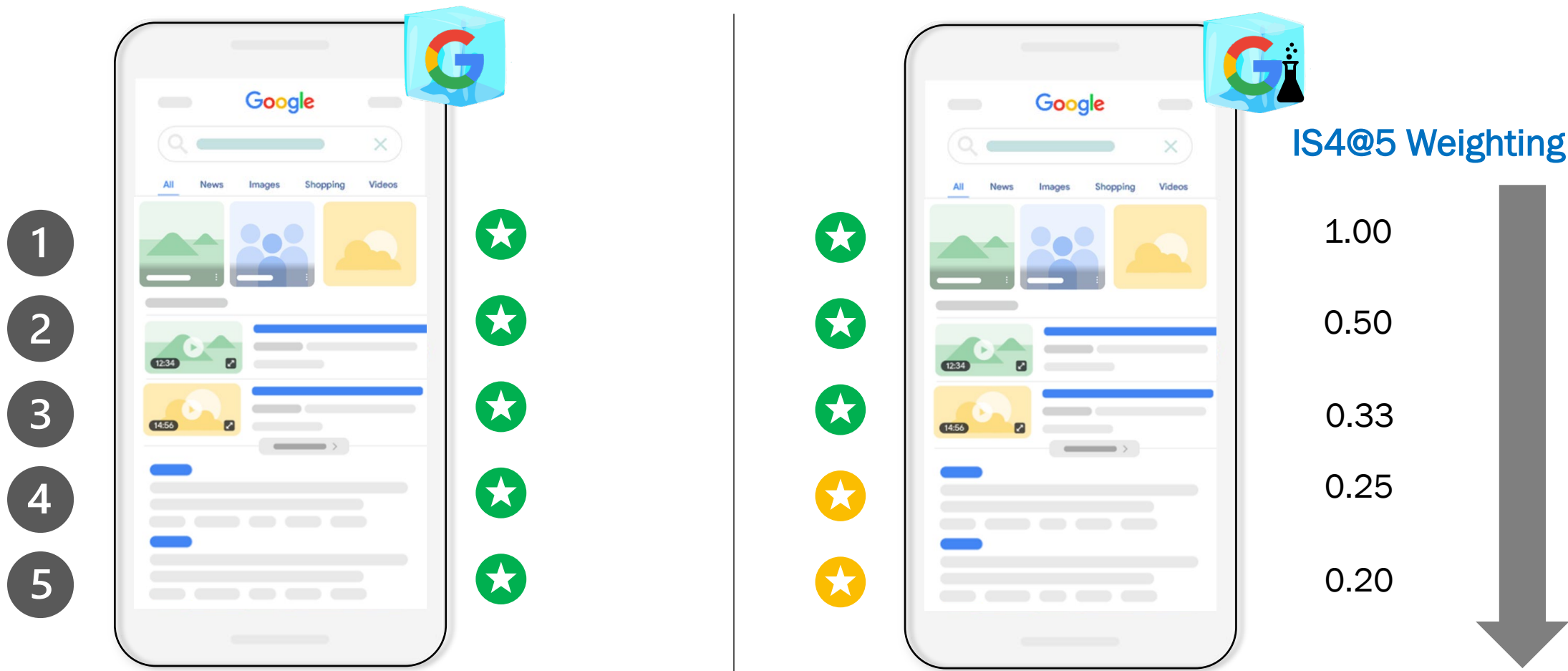


Prof. Edward
Fox
Google's Expert
Witness

“Google rates the top five positions for IS4@5 **counting both search features** like OneBoxes and **‘blue links.’**”

REDACTED FOR PUBLIC FILING

“10 Blue Links” Ranking’s Effect on IS4@5 Can Be Small



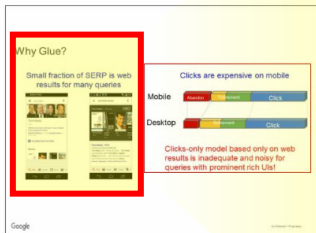
Frozen Google

Retrained Google

REDACTED FOR PUBLIC FILING

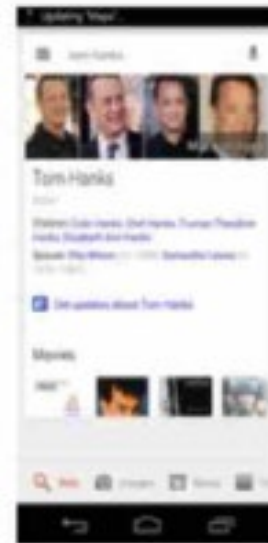
Search Features “Trigger” for Many Results

Google
Glue 2.0
glue-models@



Why Glue?

Small fraction of SERP is web results for many queries



“Small fraction of SERP is web results **for many queries**”

REDACTED FOR PUBLIC FILING

My Response to Prof. Fox's Central Conclusions

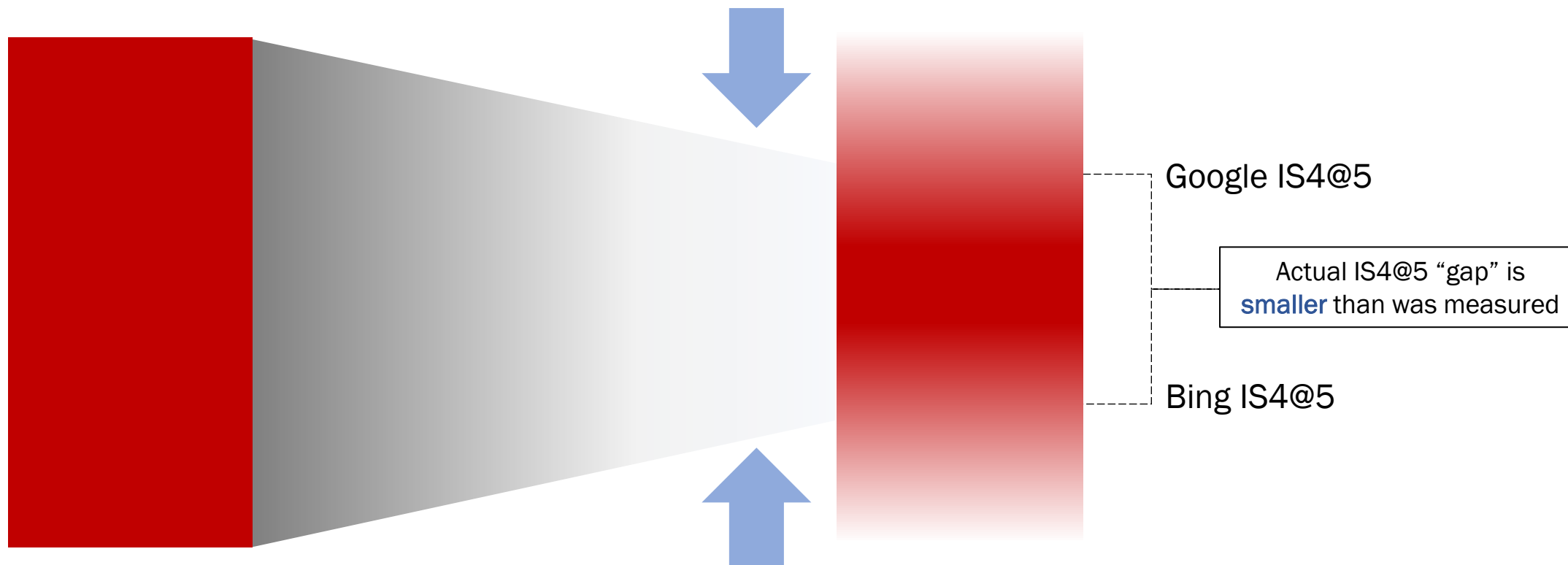
Prof. Fox's conclusions are unsupported because of:

- Unmeasured benefits of user-side data in this experiment;
- Measurement errors in the “quality gap”; and
- Important benefits of user-side data that this experiment cannot measure

REDACTED FOR PUBLIC FILING

Correcting for Measurement Errors

1. Effect of Google “teaching to the test”



2. Google's choice to "rate" **all queries** using based on **mobile presentation**

3. Google's difficulties accurately rating Bing's results

REDACTED FOR PUBLIC FILING

Teaching to the Test: Google Trains Using IS, Bing Does Not

The Six Ranking Components Retrained in the DRE

Component Name	Description	Data Used	Model Comments
Navboost	Redacted	Query-click log data is tabulated, including counts from training. Engineered functions have parameters that are learned by trying to maximize the IR ratings of the queryset result rankings they produce. A function maps each table-based Navboost signal value to a score multiplier.	Redacted
RankBrain		Trained on pairwise click preferences. Then fine-tuned on 3 rating data. Language understanding works with unigrams and bigrams.	
DeepRank		Pre-trained on document (URL, title, salient terms) data and queries/focals, and then on pairwise and pairwise click data. Then it is fine-tuned on rating 3 data. Language understanding works with word pieces.	
QBST		Trained on documents and query-click logs. Then the ranking integration is trained on rating data. Language understanding works with unigrams and bigrams.	
Term Weighting		Trained on query-click logs. Then the ranking integration is trained on rating data.	
RankEmbedBERT		Trained on documents, queries, click logs, and rating data. Makes use of salient terms and Navboost data.	

Expert Report of Edward A. Fox, App. A at Table 1

Google
DXD-26.004

Component Name	Data Used
Navboost	...Engineered functions have parameters that are learned by trying to maximize the IS ratings of the queryset result rankings they produce...
RankBrain	...Then fine-tuned on IS rating data...
DeepRank	...Then it is fine-tuned on rating IS data...
QBST	...Then the ranking integration is trained on rating data...
Term Weighting	...Then the ranking integration is trained on rating data.
RankEmbedBERT	Trained on documents, queries, click logs, and rating data...

REDACTED FOR PUBLIC FILING

Mobile Evaluation Understates Bing's Search Quality

“On Desktop, Google is comparable to Bing”

Privileged & Confidential
Search Quality (Features + Ranking): On Mobile, Google leads all Search Engines. On Desktop, Google is comparable to Bing, but leads all others

Redacted

Desktop
Mobile

Across the board,
Google outperforms
more on mobile than
desktop

Aug '20 [Search Comparative Research](#)

Google

2020

REDACTED FOR PUBLIC FILING

The Measured “Quality Gap” Does Not Account for This



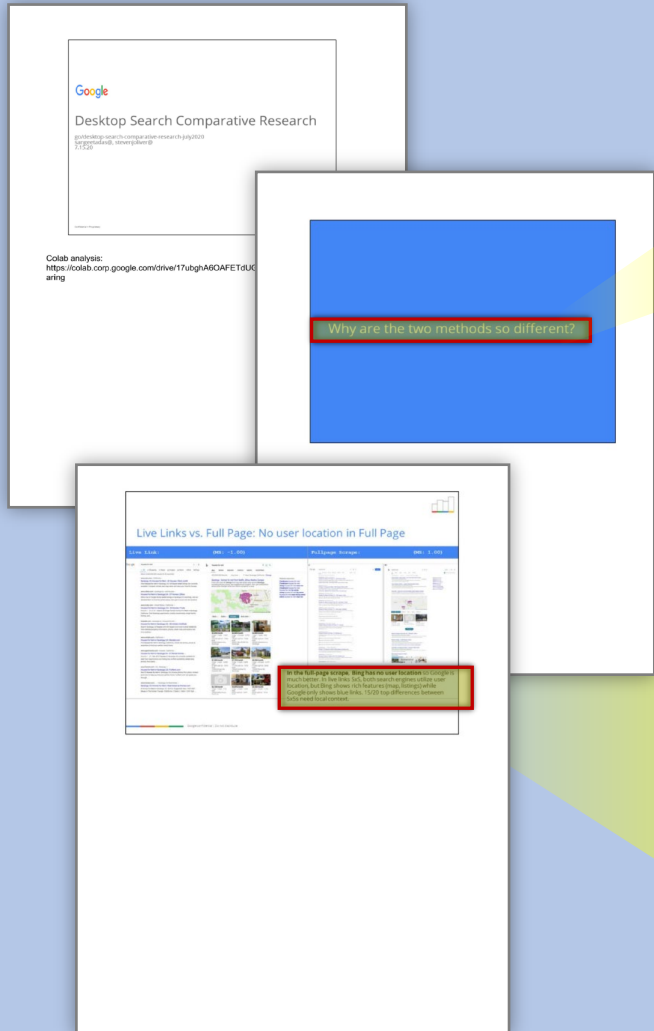
Prof. Edward
Fox
Google’s Expert
Witness

Q. You don’t know what the IS gap would be if human raters were looking at desktop presentation; right?

A. Google made a decision some years ago to do all the rater experiments with mobile. **So that’s all I know.**

REDACTED FOR PUBLIC FILING

“Scraped” Results Can Understate Bing’s Search Quality



“Why are the two methods so different?”

“In the full-page scrape, **Bing has no user location so Google is much better.**

In live links SxS, both search engines utilize user location, **but Bing shows rich features (map, listings) while Google only shows blue links.**”

2020

REDACTED FOR PUBLIC FILING

Correcting for Measurement Errors



Measured difference between Bing and Google

Correcting for measurement errors

Accounting for Unmeasured Benefits



Effect of retraining six components with less user-side data

Effect of retraining all components with less user-side data

REDACTED FOR PUBLIC FILING

My Response to Prof. Fox's Central Conclusions

Prof. Fox's conclusions are unsupported because of:

- Unmeasured benefits of user-side data in this experiment;
- Measurement errors in the “quality gap”; and
- Important benefits of user-side data that this experiment cannot measure

REDACTED FOR PUBLIC FILING

The Experiment Cannot Measure All Effects of User-Side Data

1

Effects on the Innovation Cycle

2

Effects that the IS4@5 Metric Can't Measure

3

Effects that a Frozen System Can't Measure

REDACTED FOR PUBLIC FILING

The Experiment Cannot Measure All Effects of User-Side Data

1

Effects on the Innovation Cycle

2

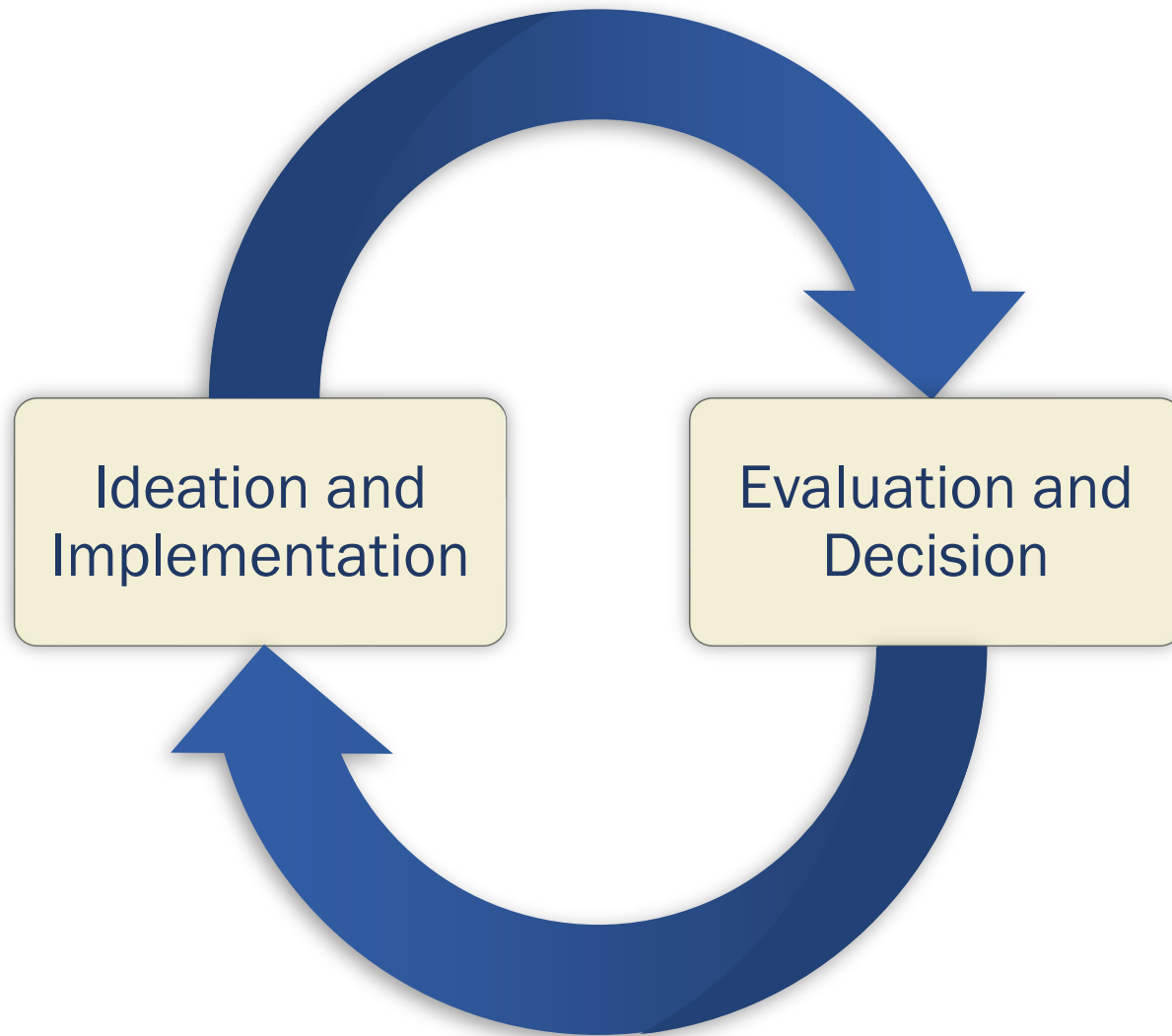
Effects that the IS4@5 Metric Can't Measure

3

Effects that a Frozen System Can't Measure

REDACTED FOR PUBLIC FILING

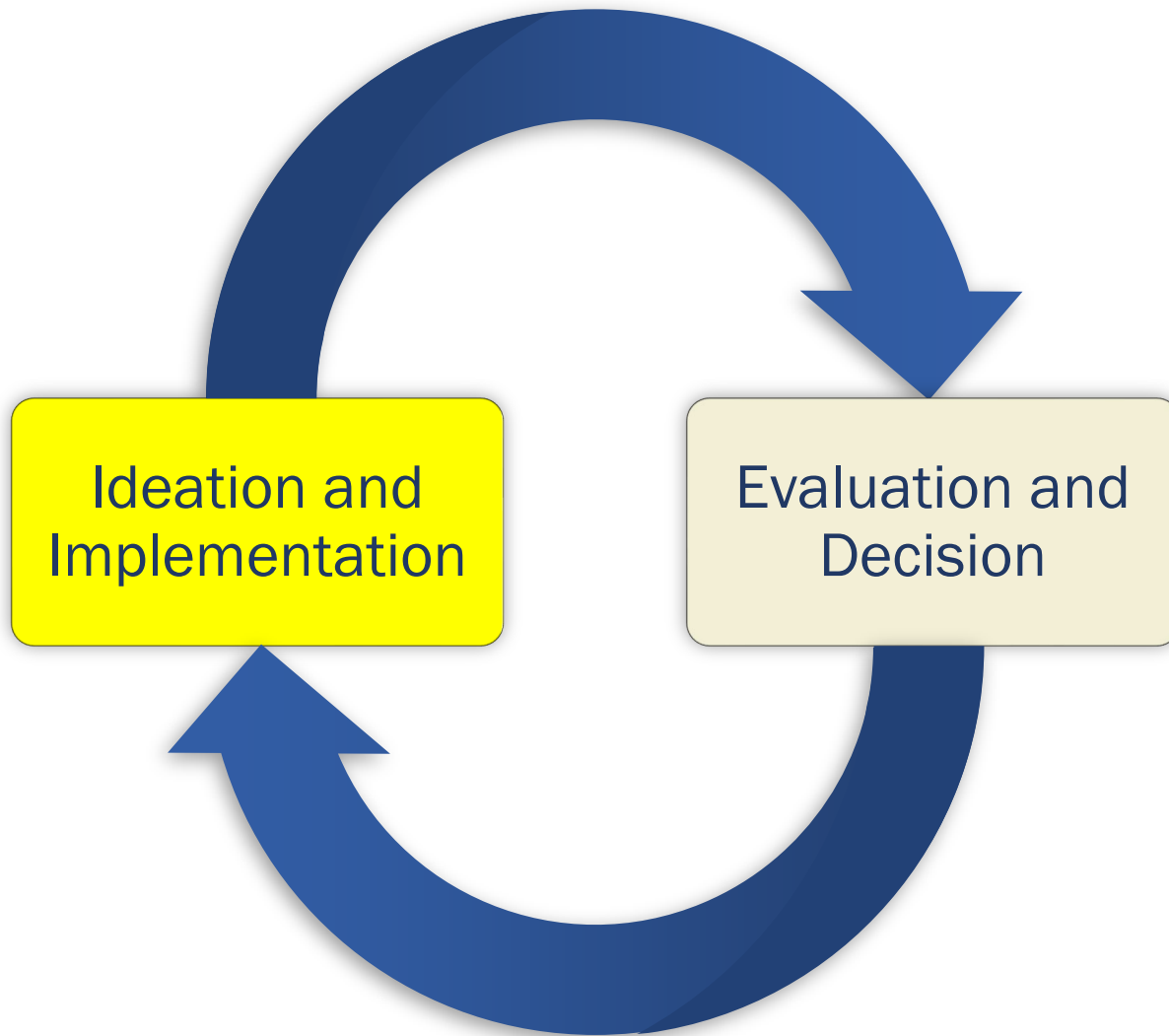
User-Side Data Benefits the Innovation Cycle



People
leverage
user-side data

REDACTED FOR PUBLIC FILING

User-Side Data Benefits the Innovation Cycle



Engineers
leverage
user-side data

REDACTED FOR PUBLIC FILING



John

Giannandrea

Apple SVP of Machine
Learning and
AI Strategy; Former
Google Head of
Search and AI



Q. ...So the **more queries** a search engine sees,
the **more opportunities it has to improve** in this
manner?

A. The **more opportunities** the engineers have **to
look for patterns and improve the algorithm,**
yeah.

REDACTED FOR PUBLIC FILING



Pandu Nayak
VP, Search



Q. ...[O]ne thing that Google might do is **look at queries for inspiration on what it might need to improve on**. Does that sound familiar?

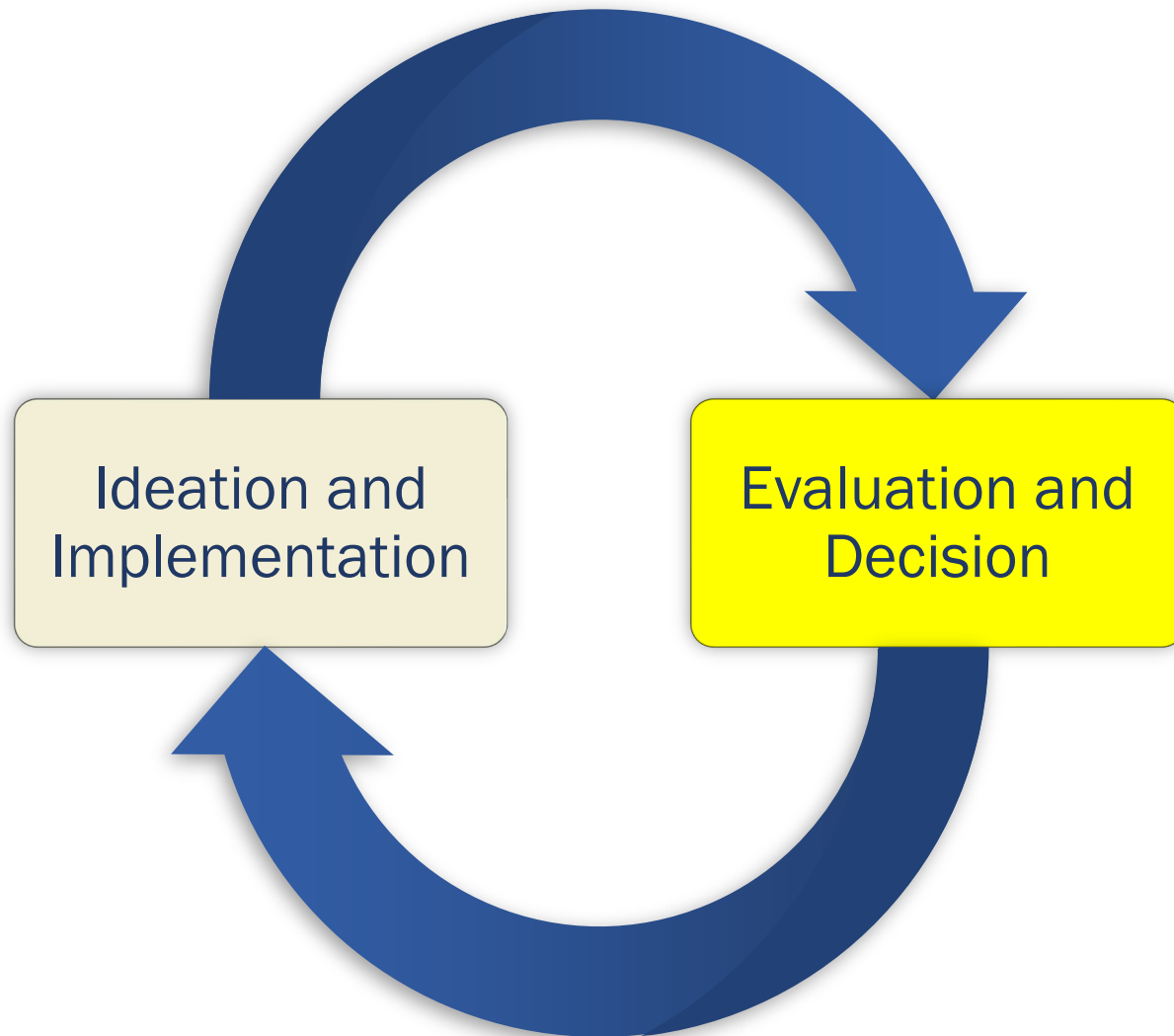
A. **Yes.**

Q. And what does that mean?

A. So we **create samples of queries** that – on which we evaluate how well we are doing overall using the IS metric, and we look at – often we look at queries that have low IS to try and understand what is going on, what are we missing here...**So that's a way of figuring out how we can improve our algorithms.**

REDACTED FOR PUBLIC FILING

User-Side Data Benefits the Innovation Cycle



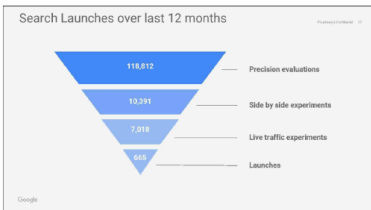
Engineers &
Managers
leverage
user-side data

REDACTED FOR PUBLIC FILING

User-Side Data Is Key to Launch Decisions

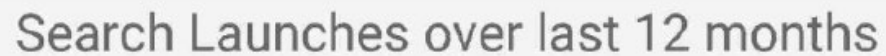


Managing Director, Global Search Ads
Google

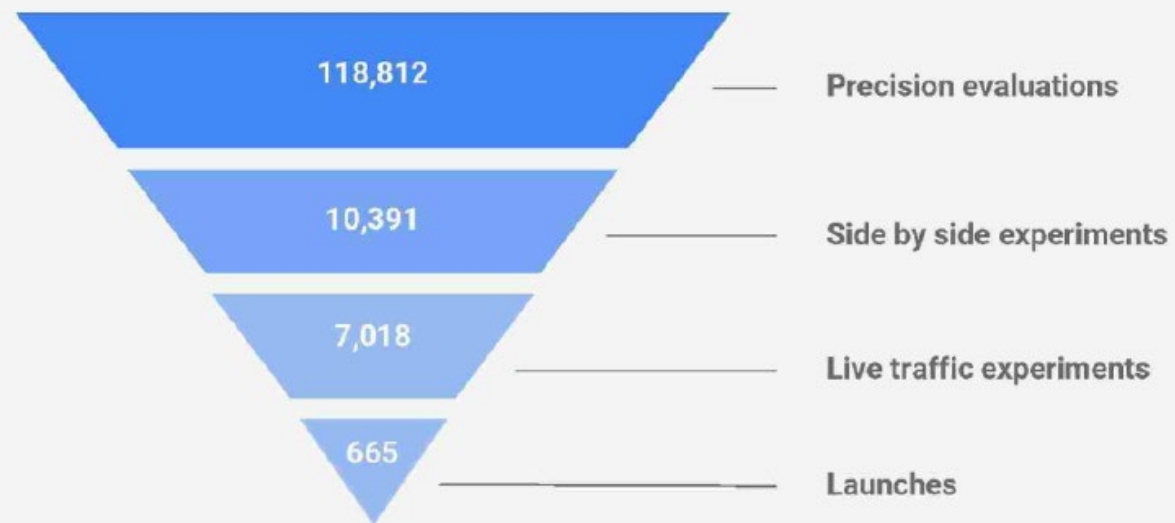


To give you a sense of the scale of this change, in a typical year, we consider over 100k changes to the search algorithm and product. Of these, 10k might seem promising and we do more detailed testing. And finally we might launch between 500 and 1k changes to the algorithm in a typical year. Each change is carefully evaluated to make sure it is an improvement to the search experience.

While I obviously can't describe to you all that we are doing, I want to give you a flavor of how we are thinking about search. We are constantly changing search with that one goal in mind - to get from the question in your mind to the information you were looking for.



Proprietary & Confidential 07



Google

2018

REDACTED FOR PUBLIC FILING

On the Role of User Interaction Data in Innovation



**Prof. Edward
Fox**
Google's Expert
Witness

Q. Whether it's innovation, better algorithms or the like you didn't study, but that's what accounts for the other 97 percent, in your view?

A. So, **I don't know what the other parts are**. I have guesses because I've worked in the field for a long time, **but it's not from user interaction data**. That's what I can tell.

REDACTED FOR PUBLIC FILING

The Experiment Cannot Measure All Effects of User-Side Data

1

Effects on the Innovation Cycle

2

Effects that the IS4@5 Metric Can't Measure

3

Effects that a Frozen System Can't Measure

REDACTED FOR PUBLIC FILING

Metrics Are Not Search Quality

IS Project Plans for 2021

...
Jan 28, 2021
Webranking Leads Meeting

IS OKR

- Launch 10 projects that deliver +0.1 IS on any one of the en-US covert set, the **18n** set, or the **longtail** set and increase IS on the en-US covert set by +1 point.

Not covering NBU OKRs for IS here.

Our Interpretation:

- IS4 is an approximation of user utility-- treat it as such.
- Aim for significant improvements when possible.

So, here's the OKR. I'm not going to read it aloud, but I'll talk about our interpretation of the OKR.

IS4 is perhaps our most important top-line ranking metric but it still is an approximation and is error-prone. We should treat it as such and always look for real user value supported by thorough analysis and other metrics.

Secondly, the 0.1 bar is meant to motivate people to work towards significant improvements when possible. At the same time, we'll continue to work on smaller, more focused improvements to address issues as they come up.

“IS4 is an approximation of user utility-- treat it as such.”

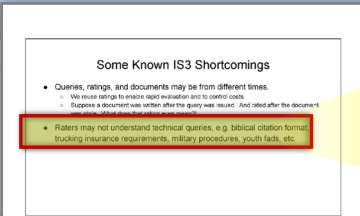
“[A]lways look for real user value supported by thorough analysis and other metrics.”

2021

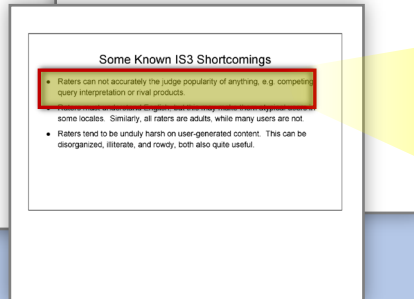
REDACTED FOR PUBLIC FILING

The IS4@5 Metric Is Only a Part of the Story

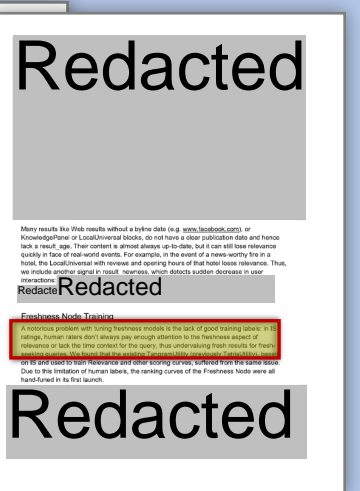
2018



“Raters may not understand technical queries”



“Raters can not accurately judge popularity of anything”



“[I]n IS ratings, **human raters don't always pay enough attention to the freshness aspect of relevance** or lack the time context for the query, thus **undervaluing fresh results for fresh-seeking queries**”

REDACTED FOR PUBLIC FILING

2021

Google Uses Many Metrics to Evaluate Search Quality

Search Quality - JG Review

June 17, 2016

Metrics

Goal is to capture user intent with metrics

Main metrics:

- *IS, PQ, Side-by-Sides, Live Experiments, Freshness*
- Use these metrics for signal development, launches, and tracking

Let's look at a few metrics...

Google

Metrics

Goal is to capture user intent with metrics

Main metrics:

- *IS, PQ, Side-by-Sides, Live Experiments, Freshness*
- Use these metrics for signal development, launches, and tracking

Let's look at a few metrics...

Google

2016

REDACTED FOR PUBLIC FILING

Live Experiment Metrics Provide Crucial Insights

Search Quality - JG Review

June 17, 2016

Live Experiments (LE)

- All Ranking experiments run LE (if possible)
- Measures position weighted long clicks
- Eval team now using *attention* as well

Redacted

Google

Live Experiments (LE)

- All Ranking experiments run LE (if possible)
- Measures position weighted long clicks
- Eval team now using *attention* as well

Redacted

Google

2016

REDACTED FOR PUBLIC FILING

The Experiment Cannot Measure All Effects of User-Side Data

1

Effects on the Innovation Cycle

2

Effects that the IS4@5 Metric Can't Measure

3

Effects that a Frozen System Can't Measure

REDACTED FOR PUBLIC FILING

Frozen Systems Are Different from Live Systems

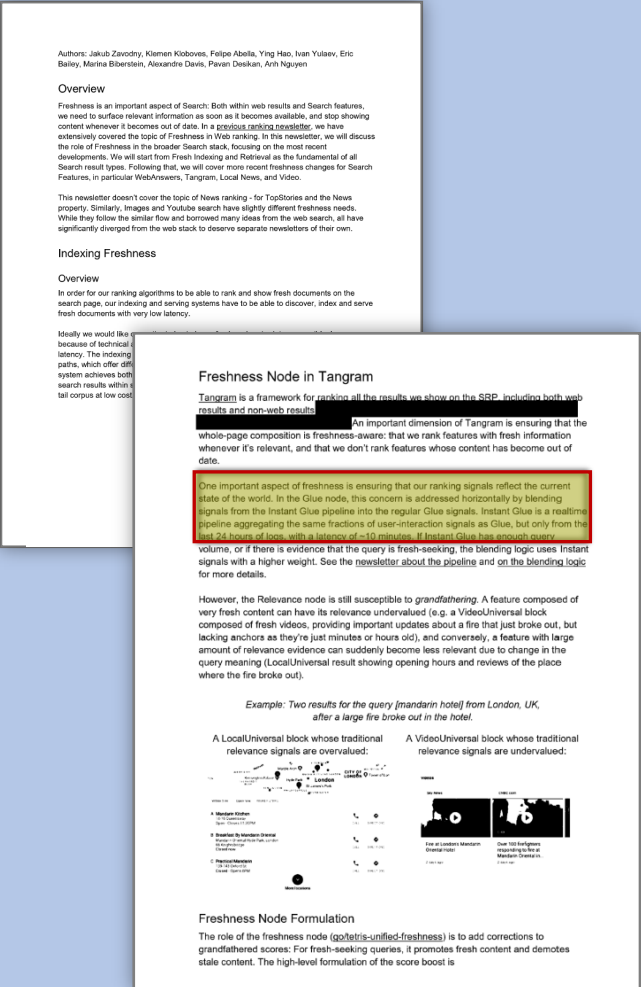
IS4@5 Difference

Redacted

Redacted

REDACTED FOR PUBLIC FILING

Frozen Systems Lack Fresh User-Side Data



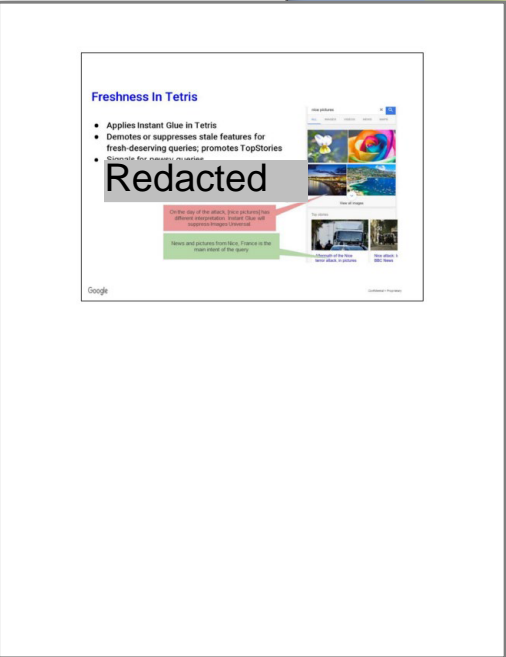
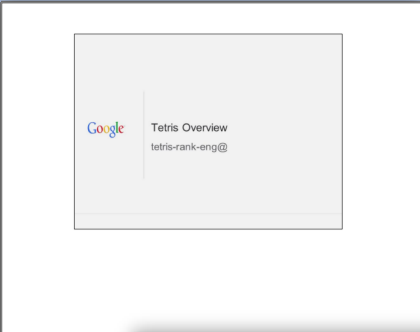
“One important aspect of freshness is ensuring that our ranking signals reflect the current state of the world.

Instant Glue is a realtime pipeline aggregating the same fractions of user-interaction signals as Glue, but only from the last 24 hours of logs, with a latency of ~10 minutes.”

2021

REDACTED FOR PUBLIC FILING

Freshness Benefits from User-Side Data



Freshness In Tetris

- Applies Instant Glue in Tetris
- Demotes or suppresses stale features for fresh-deserving queries; promotes TopStories
- Signals for newsy queries

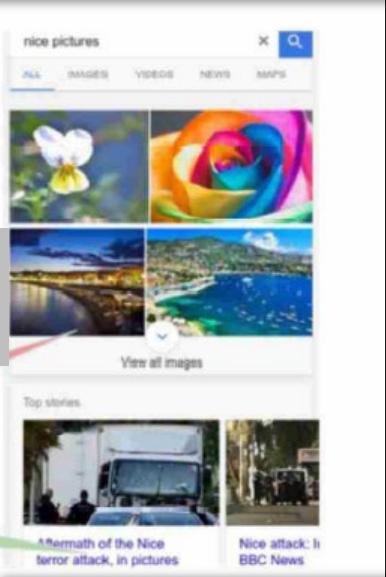
Redacted

On the day of the attack, [nice pictures] has different interpretation. Instant Glue will suppress Images Universal.

News and pictures from Nice, France is the main intent of the query

2018

On the day of the attack, [nice pictures] has different interpretation. Instant Glue will suppress Image Universal.



News and pictures from Nice, France is the main intent of the query.

REDACTED FOR PUBLIC FILING

This Experiment Can't Test Effects of User-Side Data on Freshness



**Prof. Edward
Fox**
Google's Expert
Witness

Redacted

REDACTED FOR PUBLIC FILING

Correcting for Measurement Errors



Measured difference between Bing and Google

Correcting for measurement errors

Beneficial Effects of User-Side Data this Experiment Cannot Measure

Accounting for Unmeasured Benefits



Effect of retraining six components with less user-side data

Effect of retraining all components with less user-side data

REDACTED FOR PUBLIC FILING

Prof. Fox's Third Conclusion

Vast majority of Google-Microsoft search quality gap must be explained by factors other than volume of user interaction data

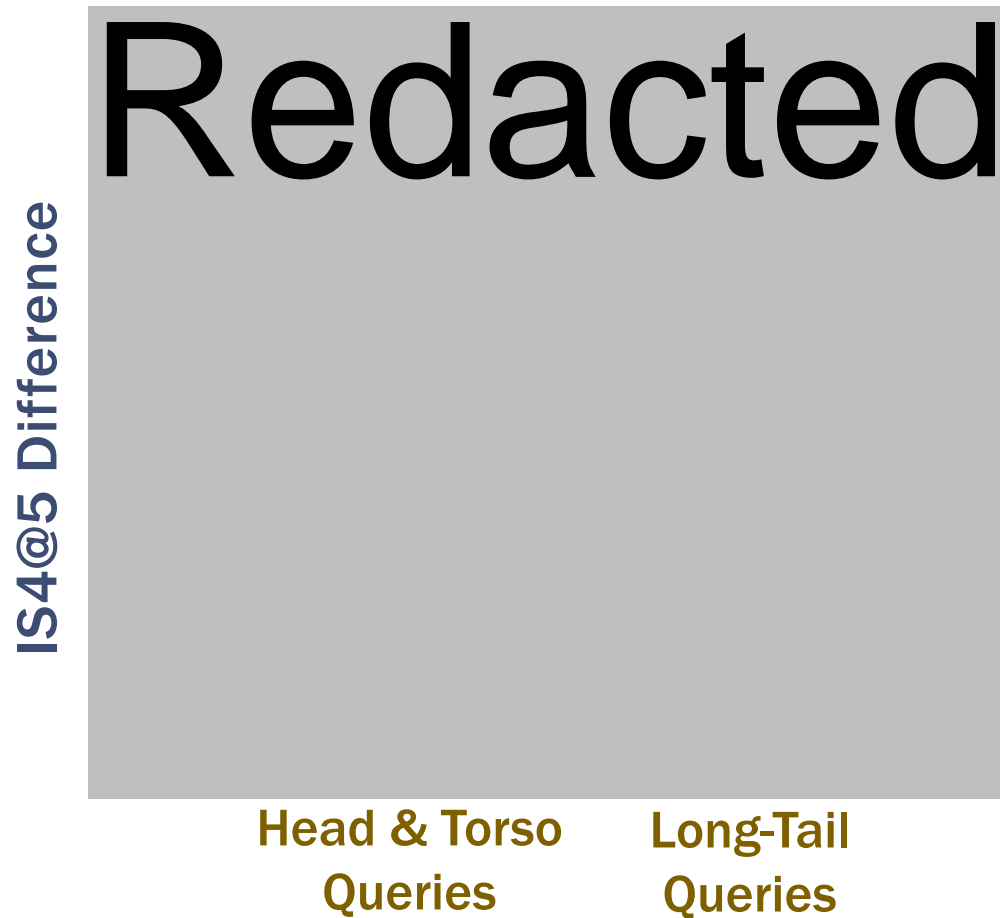
A company as efficient as Google could have search quality similar to Google even at Microsoft's scale

A company as efficient as Google but with Microsoft's scale would not meaningfully benefit from increase in user interaction data

There are diminishing returns to search quality from an increase in the quantity of user interaction data

REDACTED FOR PUBLIC FILING

The Results Show a Substantial Effect on Long-Tail Queries



Beneficial effects of user-side data can be different for different queries

REDACTED FOR PUBLIC FILING



Pandu Nayak
VP, Search



A. So we came up with the following way of thinking about it: **Wikipedia is a really important source on the web, lots of great information.** People like it a lot. **If we took Wikipedia out of our index,** completely out of our index, **then that would lead to an IS loss of roughly about a half point.** So that gives you a sense for what a point of IS is. **A half point is a pretty significant difference** if it represents the whole Wikipedia wealth of information there...

REDACTED FOR PUBLIC FILING

Prof. Fox's Final Conclusion

Vast majority of Google-Microsoft search quality gap must be explained by factors other than volume of user interaction data

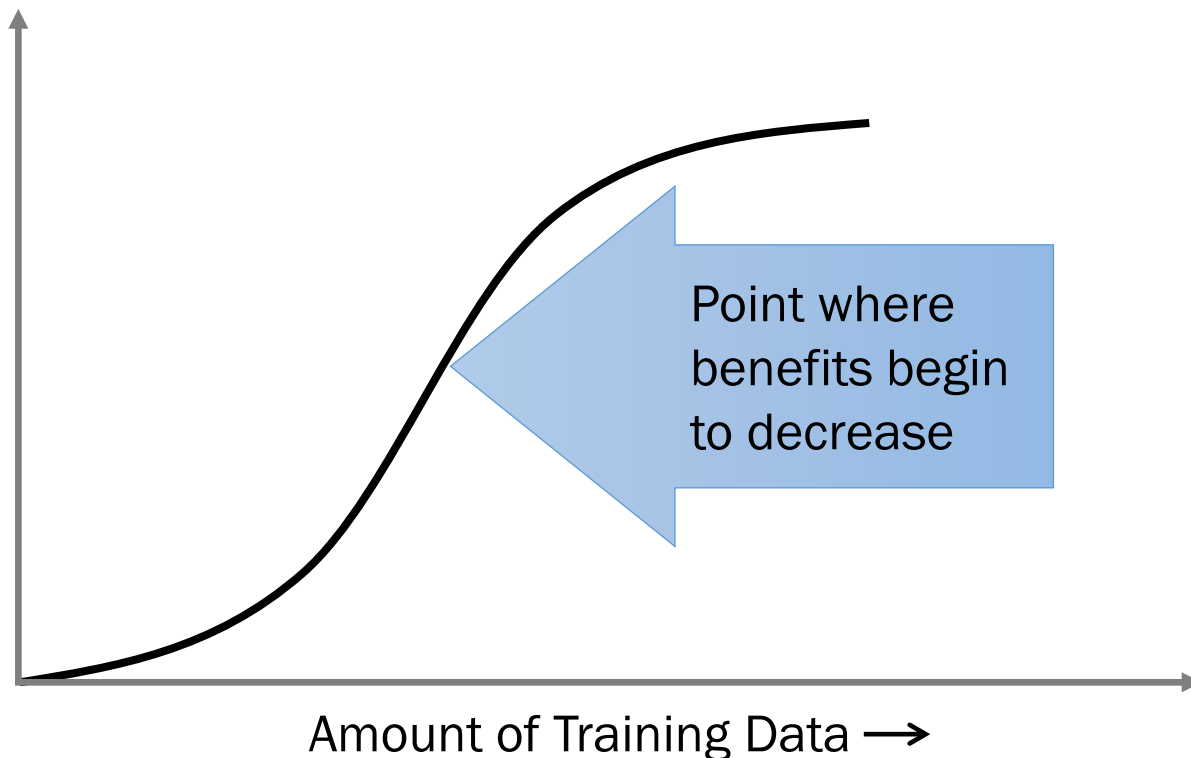
A company as efficient as Google could have search quality similar to Google even at Microsoft's scale

A company as efficient as Google but with Microsoft's scale would not meaningfully benefit from increase in user interaction data

There are diminishing returns to search quality from an increase in the quantity of user interaction data

REDACTED FOR PUBLIC FILING

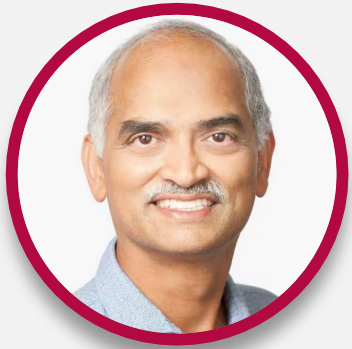
Diminishing Returns Are Not Vanishing Returns



- Benefits continue to accrue
- Benefits would be greater for tail queries, fine-grained location, etc.

REDACTED FOR PUBLIC FILING

Google's Choices Confirm Benefits Continue to Accrue



Pandu Nayak
VP, Search



Q. Google has a large collection of sessions logs. Does each click, each piece of data have the same value to Google?

A. ...And so there is this trade-off in terms of amount of data that you use, the diminishing returns of the data, and the cost of processing the data. And so usually **there's a sweet spot along the way where the value has started diminishing, the costs have gone up, and that's where you would stop.**

REDACTED FOR PUBLIC FILING

Google Clearly Gets “Returns” from User-Side Data

