# LLM-driven Relevance Engineering
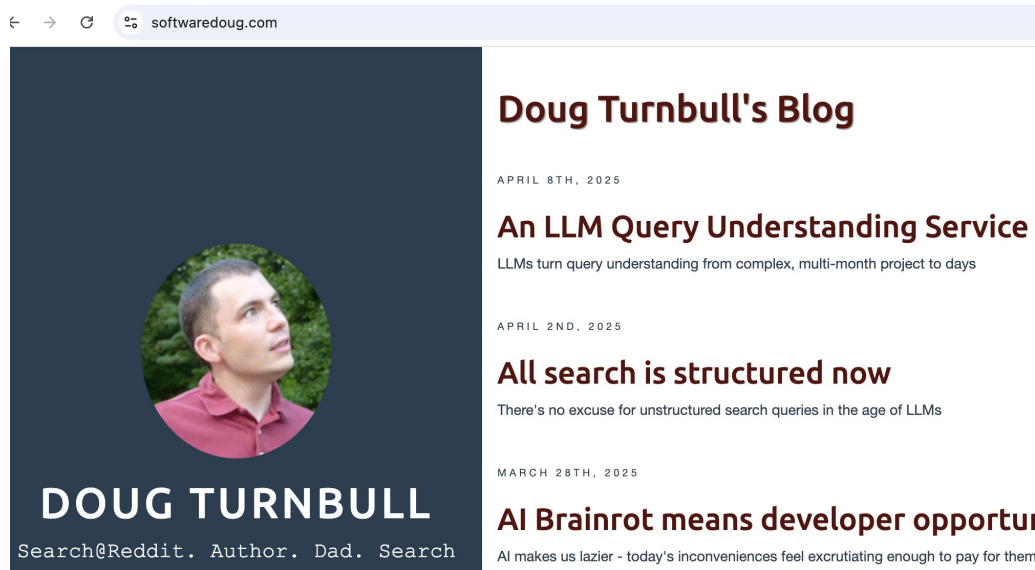
Haystack 2025

# Obligatory Bio Slide

👋 Hi I'm Doug
(@softwaredoug everywhere)
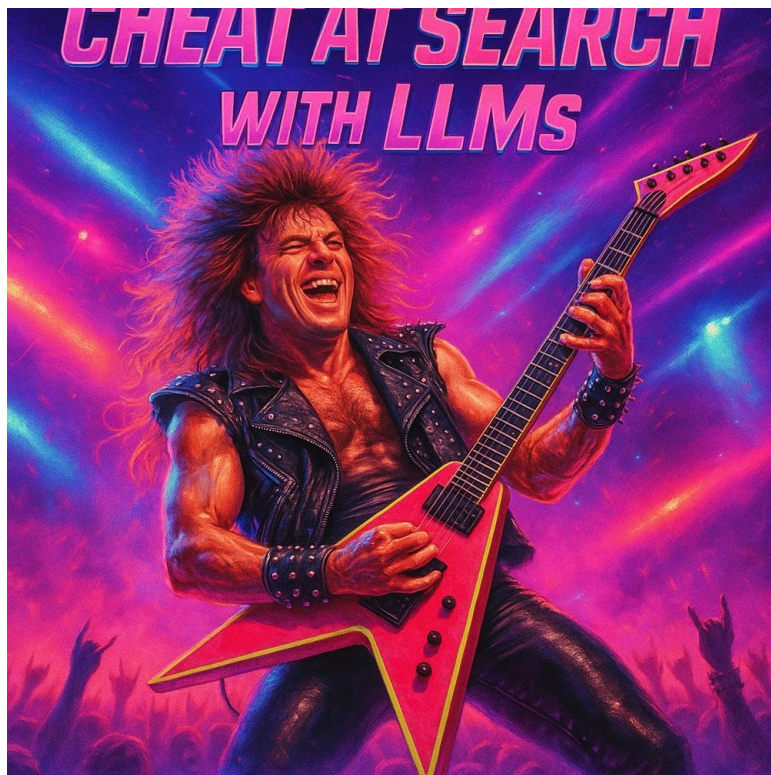
I blog here: http://softwaredoug.com

~~Live in Concert~~

# Training course:
**LLM Query understanding, content, and judging**



Obligatory Plug
https://maven.com/softwaredoug/cheat-at-search
Discount Code: **searchybird** good through Apr
(EXPIRES 1.5 weeks)

# Disclaimers, etc

- Going to talk about my specific path — what I did at Daydream (but on open dataset)

- IE not the end-all be-all of this topic *by far!*
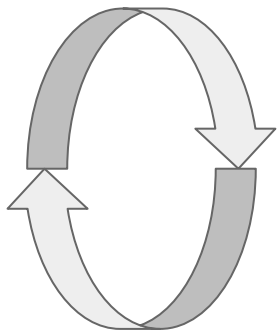
# Search labeling very time consuming

Labelers

Labelers

Labelers

| Query | Document | Rating (1-5) |
|-------|----------|--------------|
| red shoes | 👟 | 5 |
| red shoes | 👠 | 5 |
| red shoes | 🥿 | 3 |
| red shoes | 🥾 | 3 |
| red shoes | 👖 | 1 |

# New algorithm, results haven't been labeled

| Query | Document | Rating (1-5) |
|-------|----------|--------------|
| red shoes | 👟 | 5 |
| red shoes | 👞 | ?? |
| red shoes | 🩰 | ?? |
| red shoes | 🛼 | ?? |
| red shoes | 👖 | 1 |
| red shoes | 👠 | 5 |

ARE THESE GOOD RESULTS?

imgflip.com

# LLM Judge: Tighter feedback loop

LLM Judge

Instant labeling

Search Engine

🧑‍💻 Tune more

| Query | Document | Rating (1-5) |
|-------|----------|--------------|
| red shoes | 👟 | 5 |
| red shoes | 👟 | 4 |
| red shoes | 🩰 | 4 |
| red shoes | 🛼 | 2 |
| red shoes | 👖 | 1 |
| red shoes | 👠 | 5 |

# Human labeling guides LLM judge



| Query | Document | Rating |
|-------|----------|--------|
| red shoes | 👟 | 5 |
| red shoes | 👟 | 4 |
| red shoes | 🩰 | 4 |
| red shoes | 🛼 | 2 |
| red shoes | 👖 | 1 |
| red shoes | 👠 | 5 |

LLM Judge

Training Data

Actual labels used

Labeling doesn't top

# Does our judge generalize?

LLM Judge

Training Data

Is our judge any good?

| Query | Document | Human Label | LLM Label | Delta |
|-------|----------|-------------|-----------|-------|
| red shoes | 👟 | 5 | 4 | 1 |
| red shoes | 👟 | 4 | 4 | 0 |
| red shoes | 🩰 | 4 | 3 | 1 |
| red shoes | 🛼 | 2 | 2 | 0 |
| red shoes | 👖 | 1 | 3 | 2 |
| red shoes | 👠 | 5 | 4 | 1 |

# Now we have two feedback loops



Accuracy of LLM judge against labelers

Accuracy of search on generated labels

LLM Judge

Search

# Human labeling program -> LLM training program



| Query | Document | Rating (1-5) |
|-------|----------|--------------|
| red shoes | 👟 | 5 |
| red shoes | 👟 | 4 |
| red shoes | 🩰 | 4 |
| red shoes | 🛼 | 2 |
| red shoes | 👖 | 1 |
| red shoes | 👠 | 5 |

LLM Judge

Training data

Actual labels used

Encodes our domains search knowledge

# This is an _*organizational process*_ NOT magic



Understanding limits of LLM judge
as labeling source

Organizational discipline to
manage raters
(though fewer needed)

# Classic judgment list problems:

- Categorical distinctions add noise (a 1 vs 2)

- (A lot of tuning required to reduce that noise)

True of humans + LLMs!

(ESPECIALLY smaller models)

You are a search relevance judge. Evaluate the relevance of the following product based on the following criteria

5 - exact match, user would buy
4 - almost exact match, user might buy…
...

# Actual goals:

- Guidance during manual tuning of search

- Minimize cost of human labels (by doing 1-5 labels)

- Handle noisiness of human categorical labels

# Switching to pairwise LLM judge



VS

Which of these two products are more relevant for the search query "red adidas sambas"?

Product LHS Name: adidas sambas
Product LHS Description: These red adidas sambas are a fashionable fitness shoe
Product RHS Name: pink adidas sambas
Product RHS Description: Check out these hot pink shoes, you'll love them!

Agent: LHS

Pros

- Pairwise choices less noisy [1]

Downsides:

- Instead of rating N, you rate NxN

1 - "When is it better to Compare Than To Score"
https://arxiv.org/pdf/1406.6618

# Annoyingly NxN calls…

# Isn't this a lot of OpenAI calls?

vs  → OpenAI

# Enter Apple Silicon



vs



CONSIDER IT DONE

GIFWAVE.COM

Dearest laptop: which is most relevant?

RHS!

# Optimizing our Judge



Training Data

🧑‍⚖️ LLM Judge

LHS or RHS?

Turns into MANY comparisons:

HumanRating(LHS) >
HumanRating(RHS)?

Still 1-5 labels
to speed up rating

# Pairwise: An Alternative Flow

| Query | Document |
|-------|----------|
| red shoes | 👟 |
| red shoes | 👠 |
| red shoes | 🥿 |
| red shoes | 🥾 |
| red shoes | 👖 |

VS

| Query | Document |
|-------|----------|
| red shoes | 👟 |
| red shoes | 👟 |
| red shoes | 🩰 |
| red shoes | 🛼 |
| red shoes | 👖 |

# **Optimizing == What prompt to use?**

Which of these two products are more relevant for the search query "red adidas sambas"?

Product LHS Name: adidas sambas
Product LHS Description: These red adidas sambas are a fashionable fitness shoe
Product RHS Name: pink adidas sambas
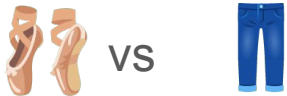Product RHS Description: Check out these hot pink shoes, you'll love them!

(What to add?)

- Product Image?
- Category?
- Reviews
- ??

More accurate decision?

(or just confusion)

# Let's try some prompts out…

(Using – **W**ayfair **AN**otated **D**ata**S**et)

| Attribute | Precision<br>(Compared to human raters. N=1000) |
|---|---|
| Just Name | 75.08% |
| Just Description | 70.31% |
| Just Category | 74.60% |
| Just Class | 70.50% |

https://github.com/softwaredoug/local-llm-judge/

# Strategies

(Using – **W**ayfair **AN**otated **D**ata**S**et)

| Attribute | Precision (Compared to human raters. N=1000) |
|-----------|-----------|
| Just Name | 75.08% |
| Just Description | 70.31% |
| Just Category | 74.60% |
| Just Class | 70.50% |
| All Fields | 78.10% |

https://github.com/softwaredoug/local-llm-judge/

# Checking twice…

Swap LHS/RHS and check for agreement

| Attribute | Precision<br>(Compared to human raters. N=1000) | Precision / Recall |
|---|---|---|
| Just Name | 75.08% / 100% | 87.99% / 58% |
| Just Description | 70.31% / 100% | 76.58% / 72.60% |
| Just Category | 74.60% / 100% | 86.1% / 69.7% |
| Just Class | 70.50% / 100% | 87.76 / 58.0% |
| All Fields | 78.10% / 100% | 91.72% / 65.2% |

# We can take it a step further

Observe what our judges are doing:

| Query/Pair | Title Judge | Desc Judge | Human Pref |
|---|---|---|---|
| Red shoes 👟 vs 👟 | LHS | Neither | LHS |
| Red shoes 🩰 VS 👟 | RHS | LHS | RHS |
| … | … | … | … |

# Train model on each judge to predict human…

| Query/Pair | Title Judge | Desc Judge | Human Pref |
|---|---|---|---|
| Red shoes  vs  | LHS | Neither | LHS |
| Red shoes  VS  | RHS | LHS | RHS |
| … | … | … | … |

Features

Predictor

Pref = f(Title_Judge, Desc_Judge, ...)

# Which of these judges has biggest impact?

(Decision tree model)

| Query/Pair | Title Judge | Desc Judge | Human Pref |
|---|---|---|---|
| Red shoes 👟 vs 👟 | LHS | Neither | LHS |
| Red shoes 🩰 VS 👟 | RHS | LHS | RHS |
| ... | ... | ... | ... |

Features

Name?

LHS

Desc?

Category?

LHS
p=0.9

RHS
p=0.51

Decision tree spits out probability of LHS/RHS

# Which of these judges has biggest impact?

| Query/Pair | Title Judge | Desc Judge | Human Pref |
|---|---|---|---|
| Red shoes 👟 vs 👟 | LHS | Neither | LHS |
| Red shoes 🩰 VS 👟 | RHS | LHS | RHS |
| … | … | … | … |

```
clf = DecisionTreeClassifier()
clf.fit(train[feature_columns],
        train['human_preference']
```

Features
(the judges)

# Which of these judges has biggest impact?

| Query/Pair | Title Judge | Desc Judge | Human Pref |
|---|---|---|---|
| Red shoes 👟 vs 👟 | LHS | Neither | LHS |
| Red shoes 🩰 VS 👟 | RHS | LHS | RHS |
| … | … | … | … |

Trying to predict

```
clf = DecisionTreeClassifier()
clf.fit(train[feature_columns],
        train['human_preference']
```
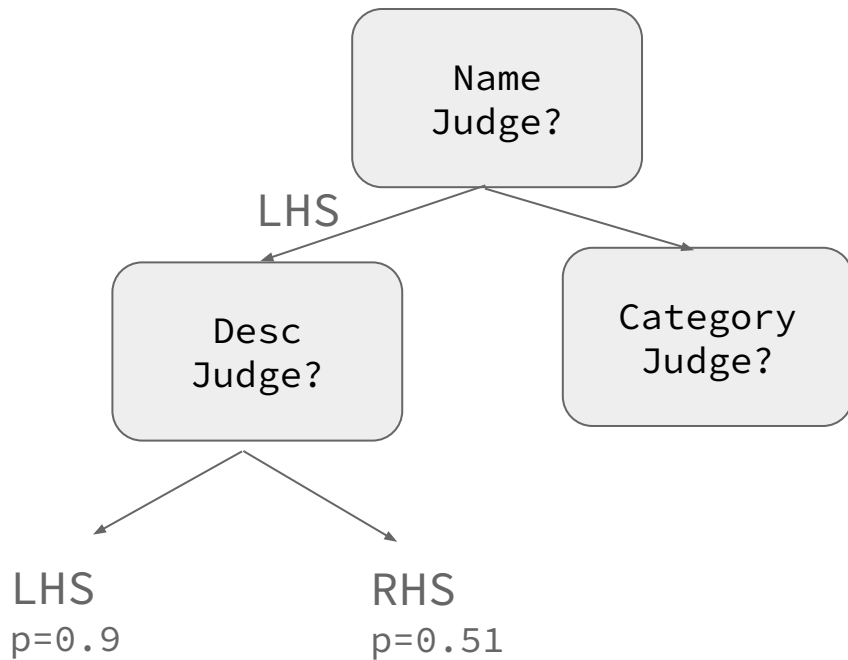
# Adding decision tree model

| Attribute | Checking Once (Compared to human raters. N=1000) | Checking Twice Precision / Recall |
|---|---|---|
| Just Name | 75.08% / 100% | 87.99% / 58% |
| Just Description | 70.31% / 100% | 76.58% / 72.60% |
| Just Category | 74.60% / 100% | 86.1% / 69.7% |
| Just Class | 70.50% / 100% | 87.76 / 58.0% |
| All Fields | 78.10% / 100% | 91.72% / 65.2% |
| Decision Tree | | 94.10% / 44.9% |

Trained on above, p > 0.9 from decision tree

Accuracy isn't
the point

# Feature Insight…

# Feature Insight…



Name
Judge?

LHS

Desc
Judge?

Category
Judge?

These are the original
features

LHS
p=0.9

RHS
p=0.51

# A lot of relevance work:

Current index:

```
{

    "Name": "Adidas Sambas"
    "Description": "These
cool red shoes…"
    "Category":
Footwear/blah/blah
    "Class": "Athletic
Footwear"

}
```
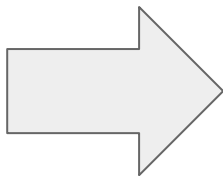
...10% of our queries mentioned color...

# LLM Judge(s) have blind spot with color

Current index:

```
{

    "Name": "Adidas Sambas"
    "Description": "These
cool red shoes…"
    "Category":
Footwear/blah/blah
    "Class": "Athletic
Footwear"

}
```
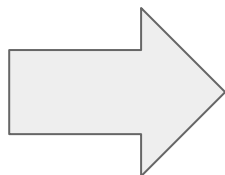
Eval for
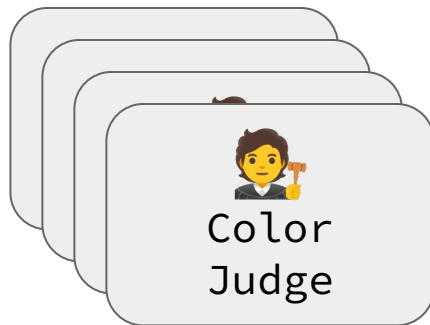*Red adidas sambas*

LLM
Judge(s)

No color to eval!

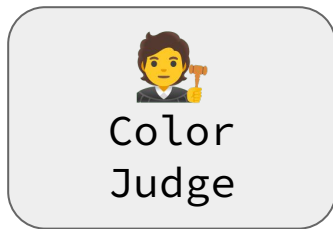# But if we added color to the index…

Current index:

```
{

    "Name": "Adidas Sambas"
    "Description": "These
cool red shoes…",
    "Category":
Footwear/blah/blah,
    "Class": "Athletic
Footwear",
    "Color": "Red"

}
```

Eval for
*Red adidas sambas*
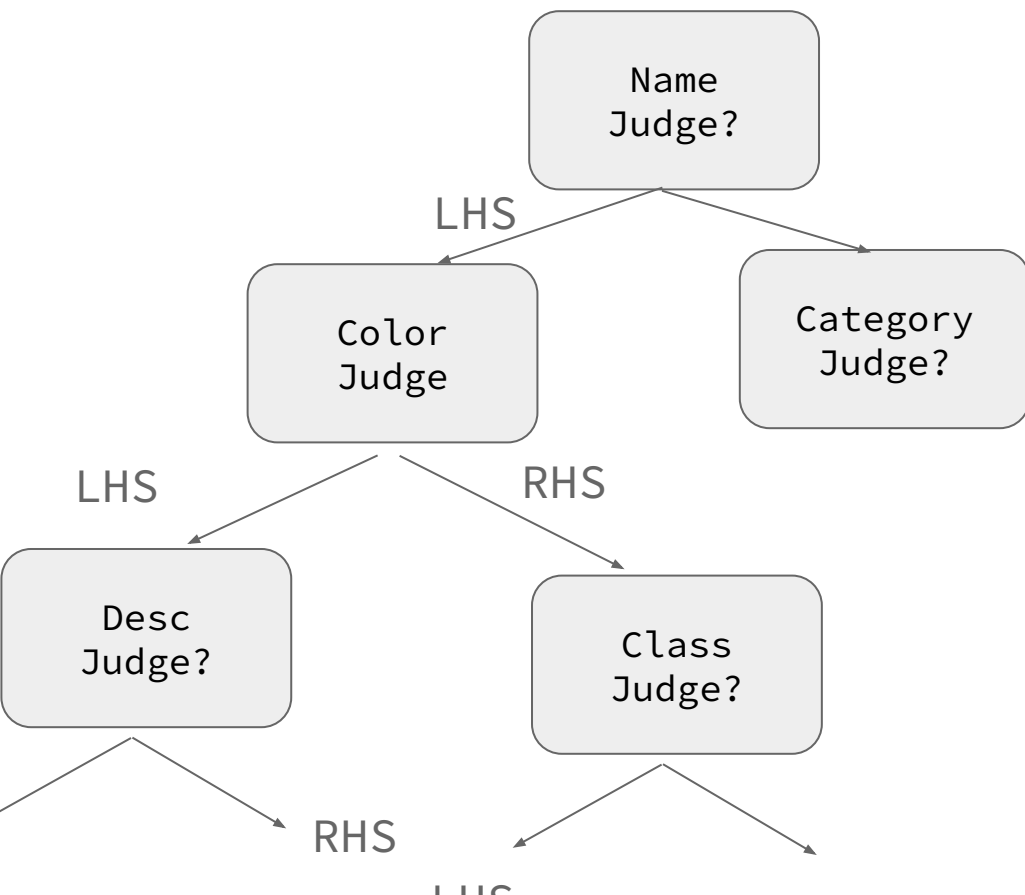


Color Judge

# With a dumb judge just to do…

Color
Judge

Which of these is more
relevant for query: red
adidas sambas

LHS Color: red
RHS Color: brown

# We could see its impact...



Name
Judge?

LHS

Color
Judge

Category
Judge?

LHS

RHS

Desc
Judge?

Class
Judge?

RHS

LHS

This model:

- Better predicts human labels by 2%

- Show color has outsized impact on decision tree eval

# Adding decision tree model

| Attribute | Checking Once (Compared to human raters. N=1000) | Checking Twice Precision / Recall |
|---|---|---|
| Just Name | 75.08% | 87.99% / 58% |
| Just Description | 70.31% | 76.58% / 72.60% |
| Just Category | 74.60% | 86.1% / 69.7% |
| Just Class | 70.50% | 87.76 / 58.0% |
| All Fields | 78.10% | 91.72% / 65.2% |
| Just Color | | 65% / 25% |
| Decision Tree + **Color** | | 95.70% / 75.9% |

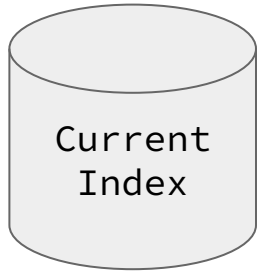Add Color to decision tree explains more of the problem

*made up stats

# We haven't done the "hard work" of search

- How will we use color if it's in the index?

- How will we balance it against other factors?

- How do we map colors in query to colors in document? Query understanding? Extraction? Etc etc?
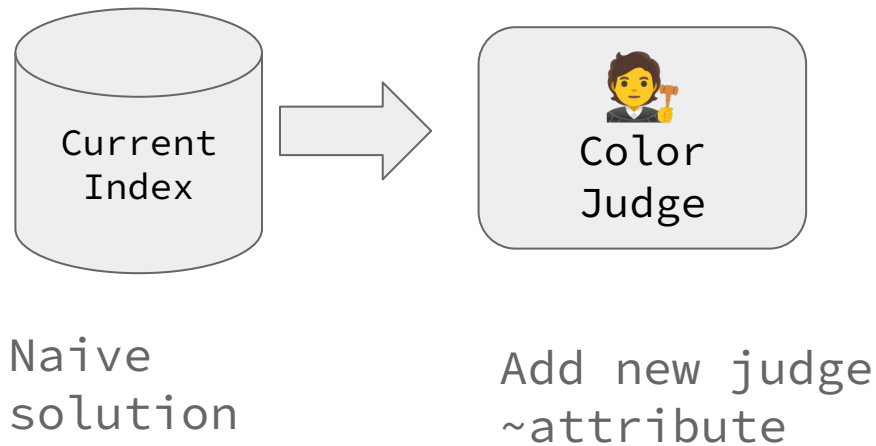
# LLM Driven Relevance Engineering
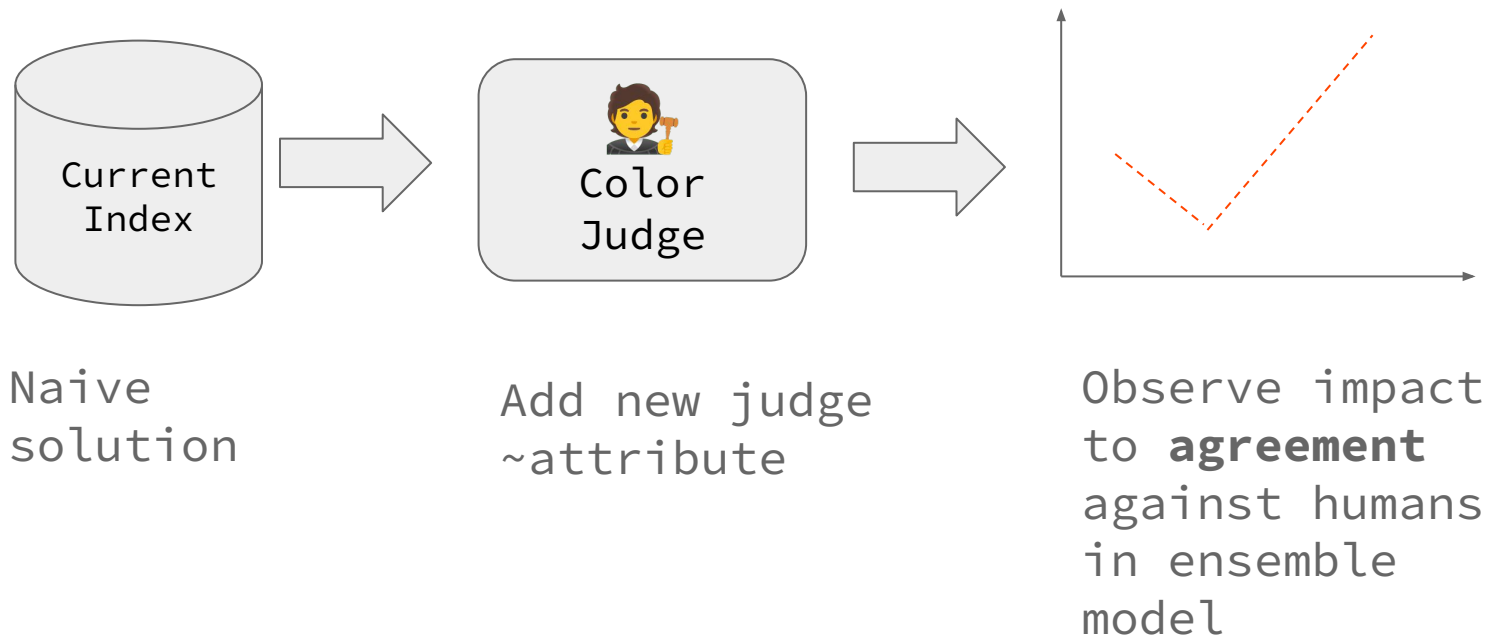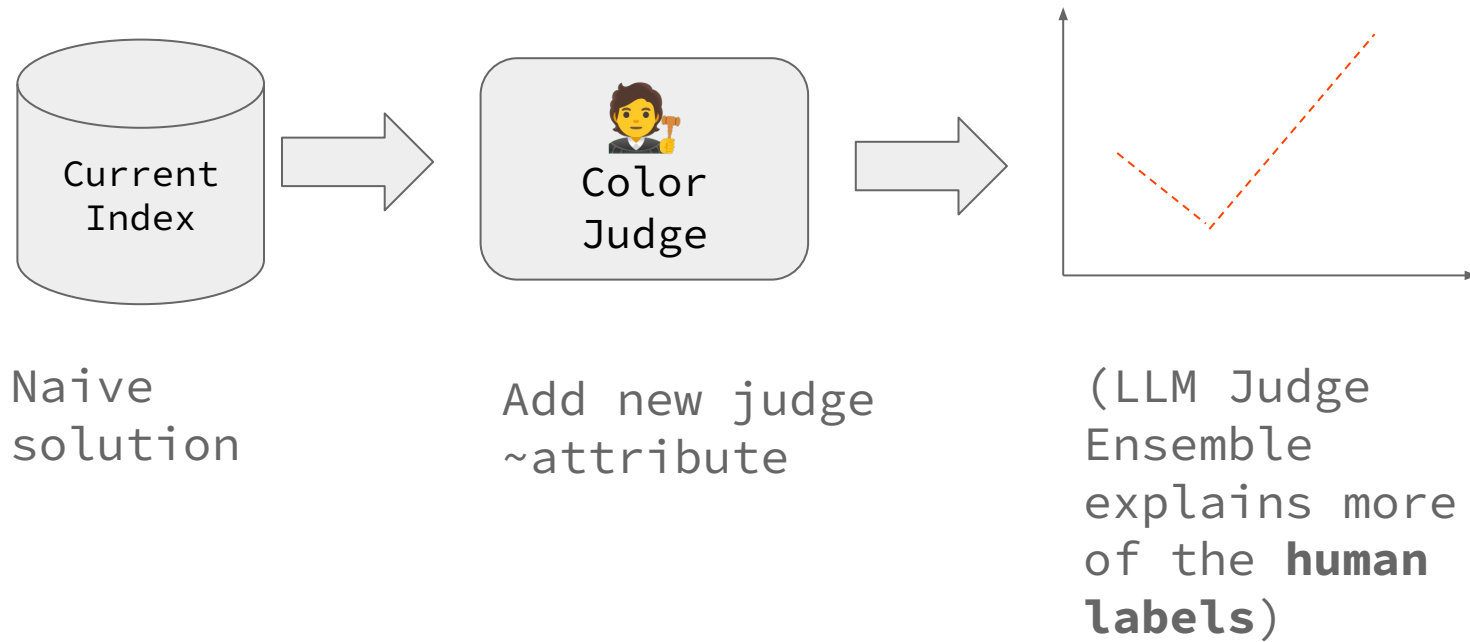
# Leading to LLM Driven Relevance Engineering

Current
Index

Naive search
solution

# Leading to LLM Driven Relevance Engineering
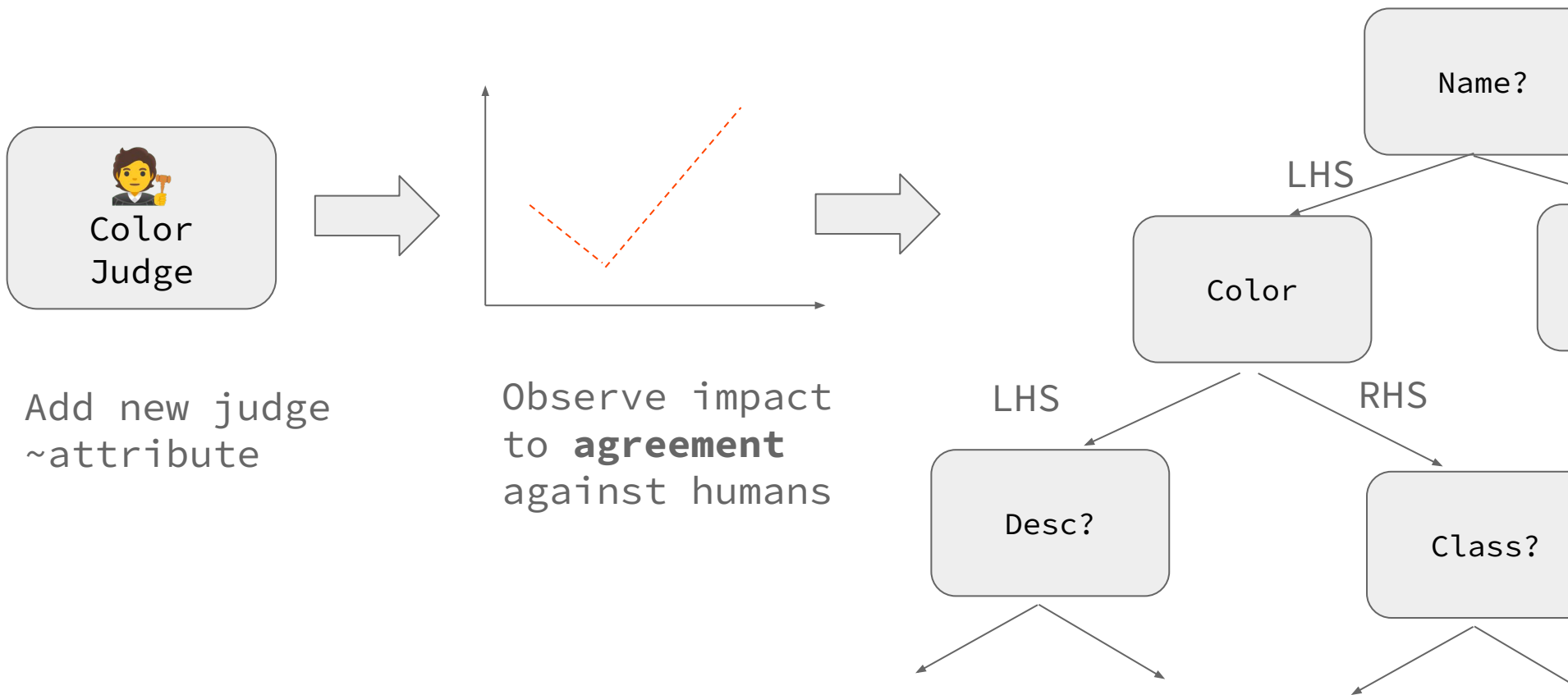
Current
Index

Color
Judge

Naive
solution

Add new judge
~attribute

# Leading to LLM Driven Relevance Engineering



Naive
solution

Add new judge
~attribute

Observe impact
to **agreement**
against humans
in ensemble
model

# Leading to LLM Driven Relevance Engineering



Naive
solution

Add new judge
~attribute

(LLM Judge
Ensemble
explains more
of the **human
labels**)

# Finally seeing problem structure…



Color
Judge

Add new judge
~attribute

Observe impact
to **agreement**
against humans

Name?

LHS

Color

LHS                    RHS

Desc?                  Class?

# To the next thing



**Recency**
Judge

Add new judge
~attribute

Observe impact
to **agreement**
against humans

Name?

LHS

Color

LHS

Desc?

RHS

**RECENCY**?