



# Using Spark-Solr at Scale: Productionizing Spark for Search with Apache Solr

FAIZAN AHMED



# About Me

## Role

Data Engineer on Flipp's Search Team

## Education

M.Sc in Computer Science from University of Waterloo

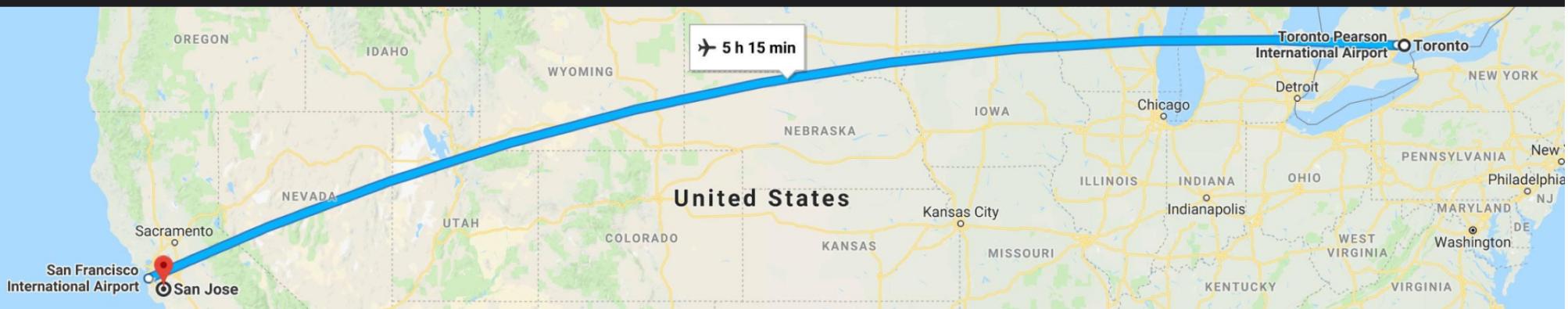
**M** [medium.com/@faizanahmed\\_18678](https://medium.com/@faizanahmed_18678)

## Responsibilities

Building pipelines to support machine learning models

Developing scalable systems and working on products to enhance search relevancy.

Core work includes using Spark for machine learning at scale & data warehousing for analytics.



# Agenda

---

- About Flipp
- The Search Experience
- The Problem Domain
- The Spark-Solr Journey
  - Solution Architecture
  - Deep Dive: Content Categorization
  - Deep Dive: Understanding User Intent
- Improvements
- Lessons Learned
- Key Takeaways

# About Flipp



## Digital Shopping Marketplace

The Flipp App is a digital shopping marketplace with over 40M+ downloads.



## Storefront Visual Merchandising

Enterprise Platform converts analog content to provide a mobile browsing experience, which showcases retailers' digital storefront.



## The Flipp Platform

Flipp's platform is used by over 90% of tier one retailers in North America.

# Search Platform At a Glance

- 2M Searches/day
- 30M Indexed Documents
- Team of 5 engineers

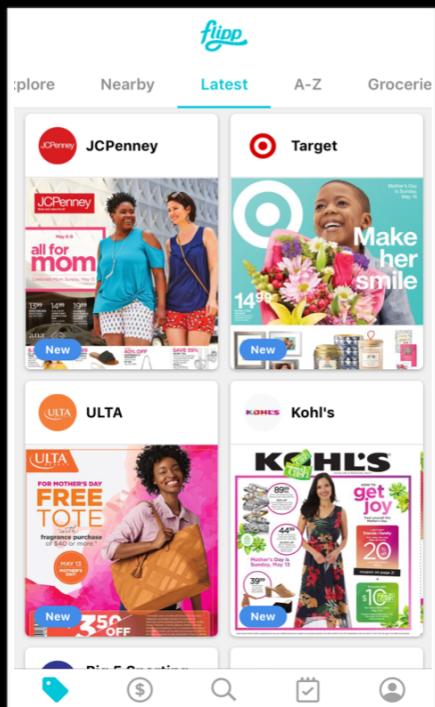


# Searchable Entities

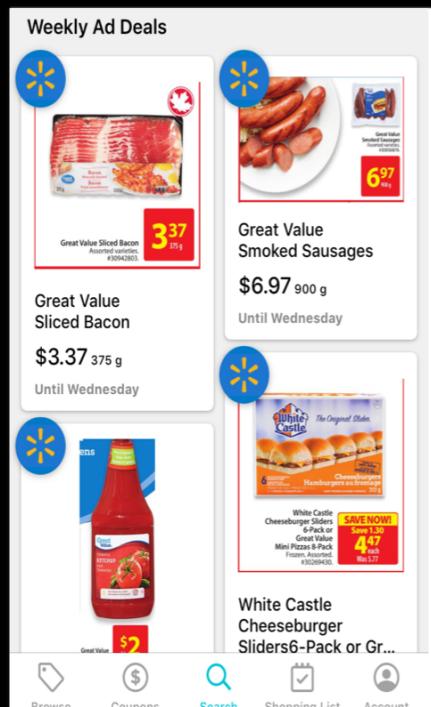


# The Search Experience

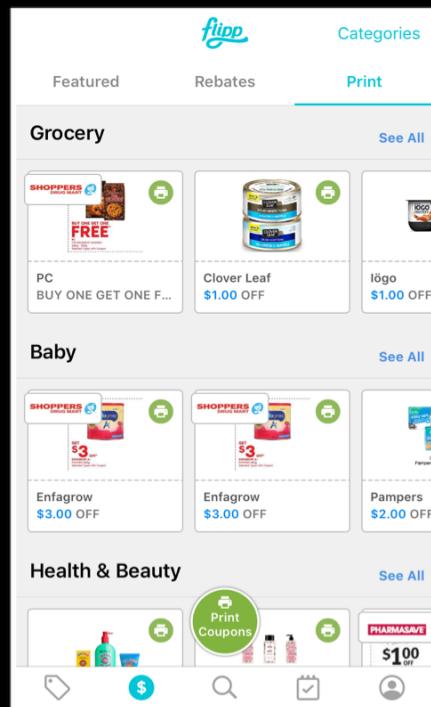
## Merchant Flyers



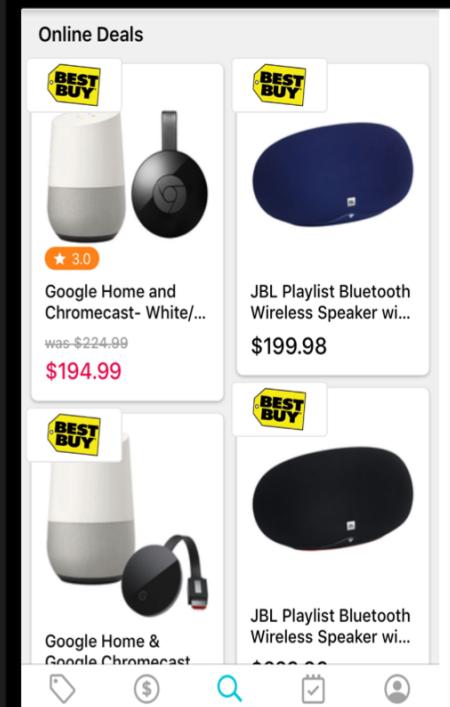
## Items



## Coupons

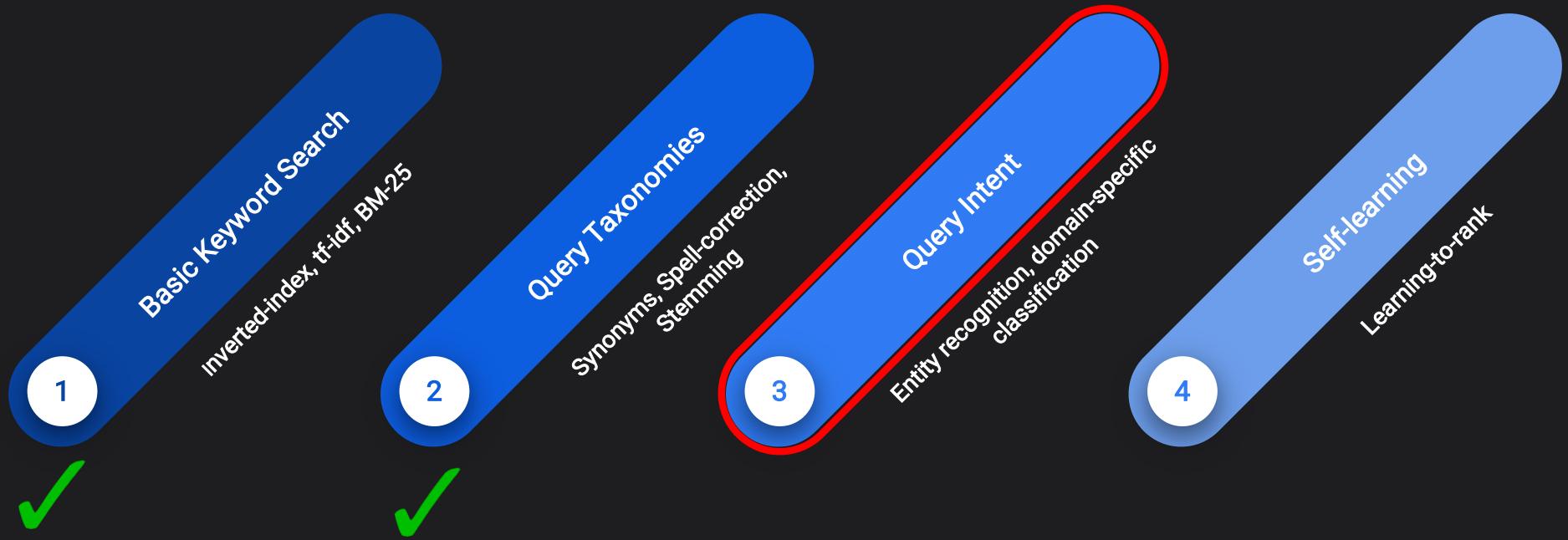


## Online Deals



# The Problem Domain

# Search Platform Evolution



Discovering and leveraging query intent is dependent on domain specific **user & content signals**.

# Signals: The Building Blocks Of Search

User Signals



Content Signals



Relevant Search

# Content Signal Challenges: Product Categorization

- Non-uniform *category* structure

Every item would have a *merchant-specific* category



Consumer Electronics -> Consumer Audio Headphones

Audio-> Earbuds & In-Ear Headphones-> Headphones

“*Bose SoundSport In-Ear Sport Headphones*”

# Content Signal Challenges: Product Categorization

- Content tagged with correct categories lacked complete hierarchy



Tagged as "Kale"



Tagged as "Greens"

Categorization hierarchy for "Kale"

Food, Beverages & Tobacco > Food Items > Fruits & Vegetables > Fresh & Frozen Vegetables > Greens > Kale

# User Signal Challenges

- Unable to leverage Solr features



```
{!boost b=log(popularity)} foo
```

“cheese” -> Dairy Products -> Cheese

“shredded cheese” -> ?

- Query categorization not deep enough



Food, Beverages & Tobacco > Food Items > Fruits & Vegetables



# What Was Needed?



ML, ETL & Analytics



Combine Signals



Spark-Solr Integration



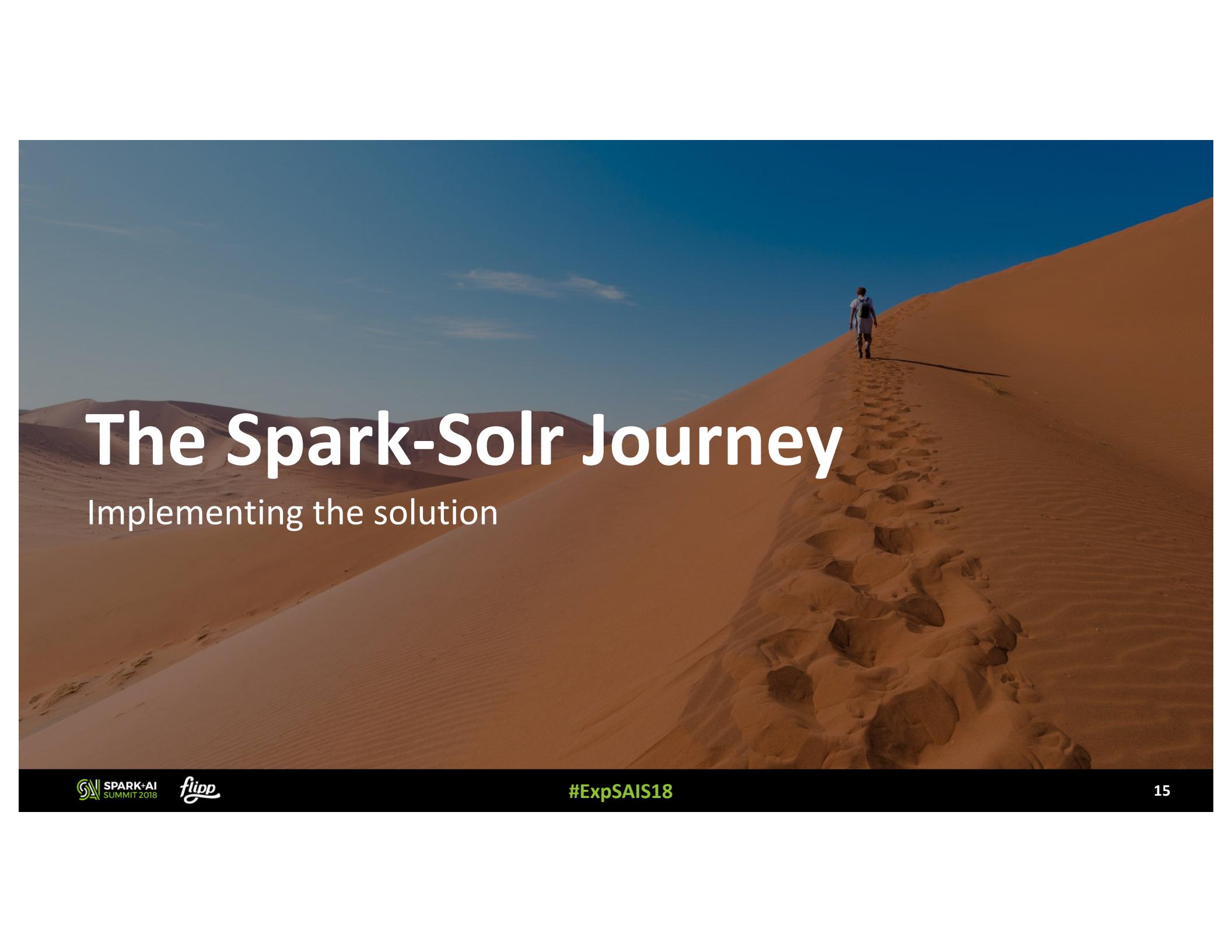
Reduce Productionization Time



Scalable & Efficient



Cost-friendly

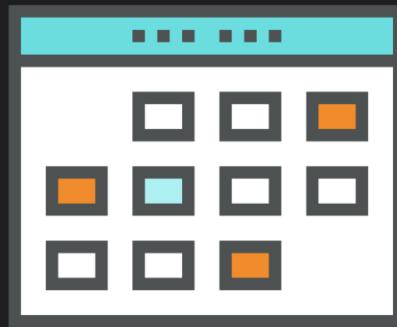
A photograph of a vast desert landscape with large, smooth sand dunes under a clear blue sky. A single person is seen from behind, walking up the side of a prominent sand dune, leaving a trail of footprints in the sand.

# The Spark-Solr Journey

Implementing the solution

# Action Plan

Content



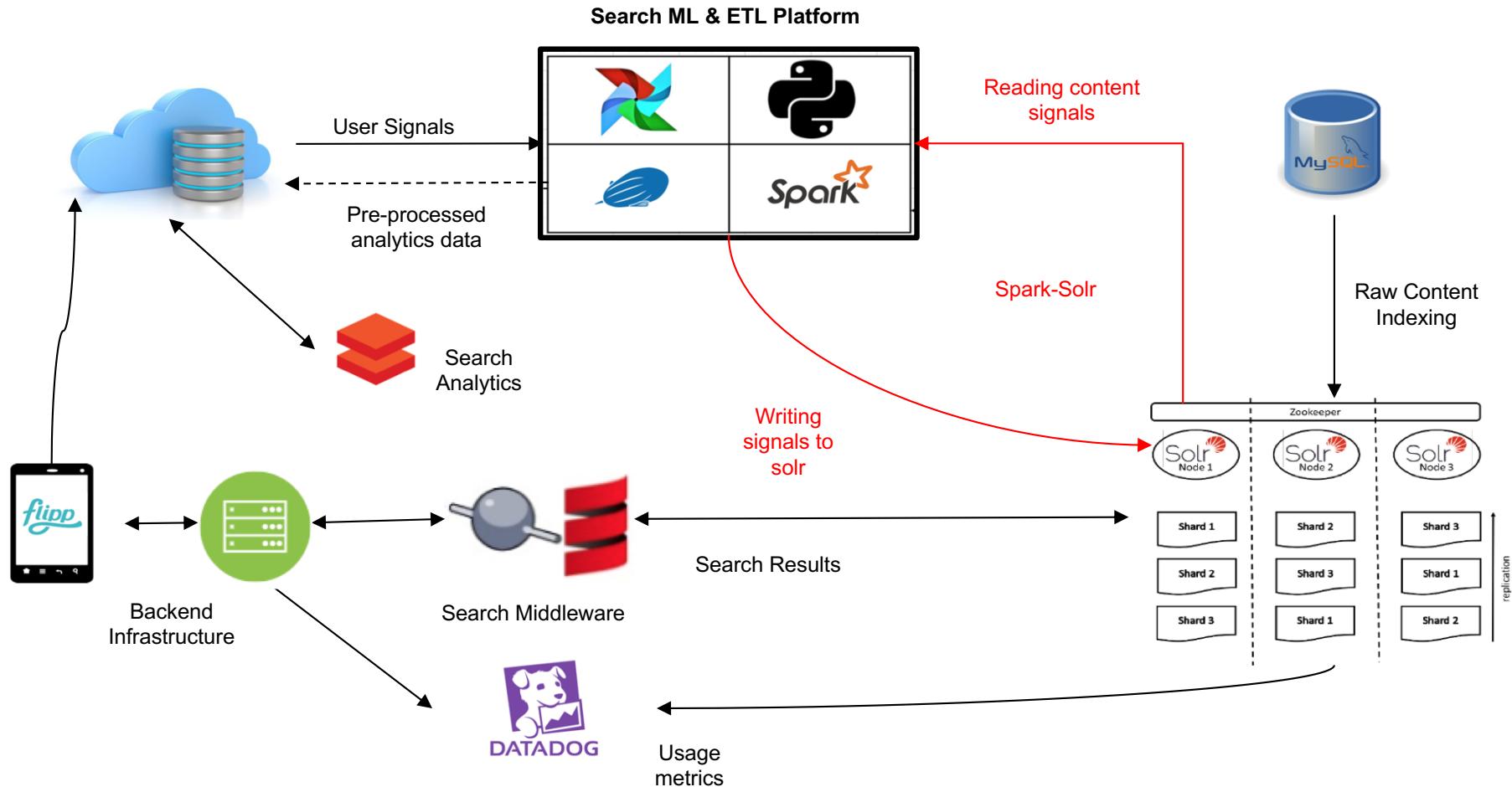
Context

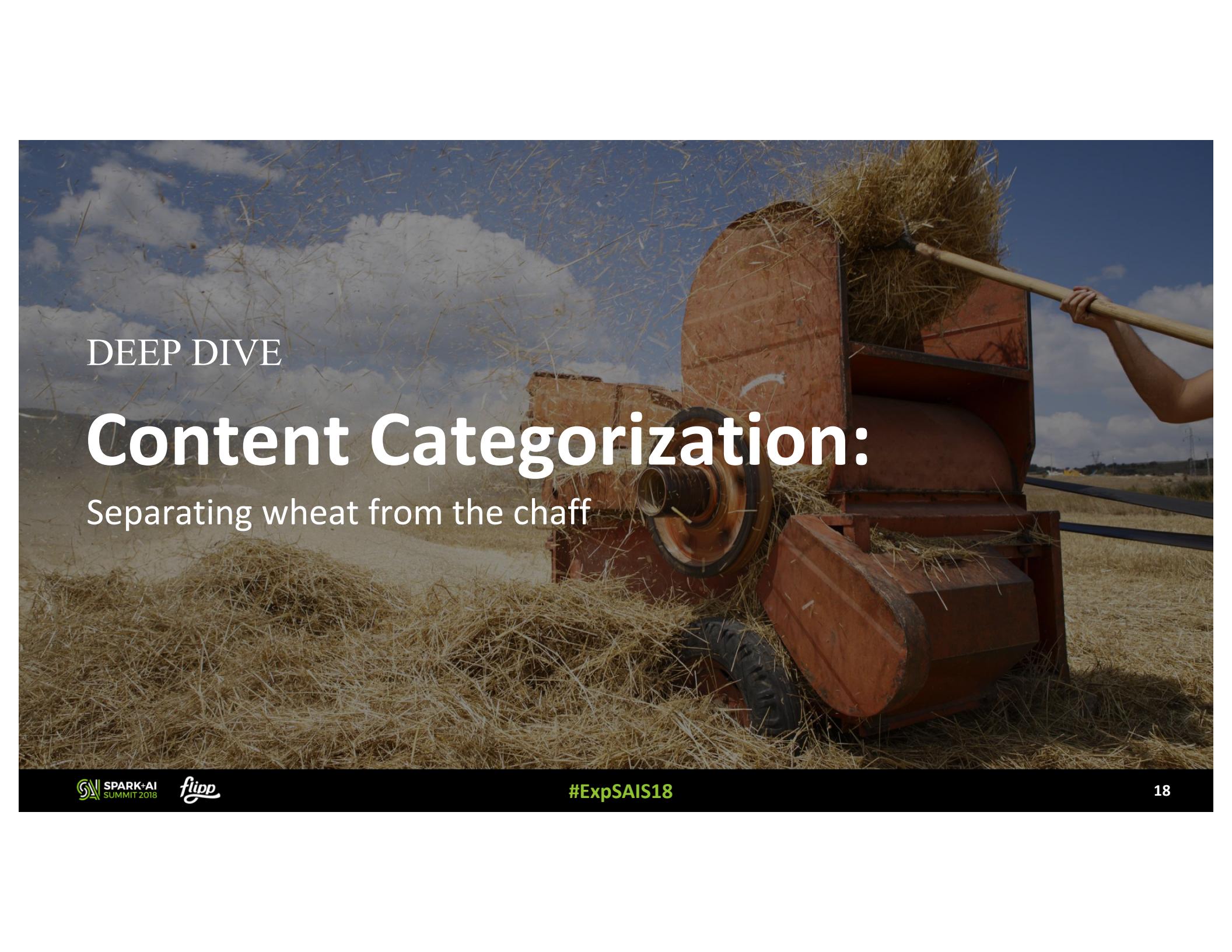


Collaboration



# The Spark-Solr Architecture



A photograph of a person using a long wooden pitchfork to move a large pile of straw or chaff from a wooden cart onto a larger pile in a field under a blue sky with white clouds.

DEEP DIVE

# Content Categorization:

Separating wheat from the chaff

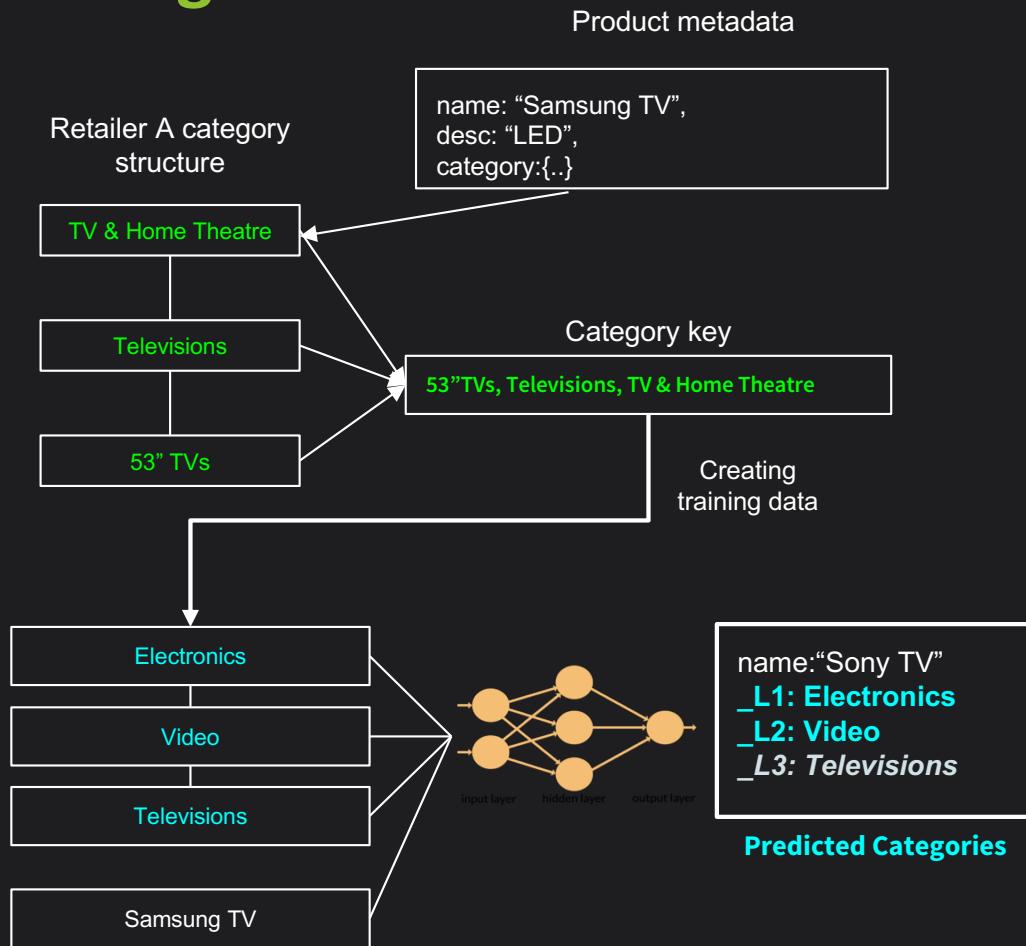
# What Raw Content In Solr Looks Like

```
"docs": [
  {
    "name_v": [
      "Samsung The Frame 55\" 4K UHD HDR LED Tizen Smart TV (UN55LS003AFXZC) - Charcoal Black"
    ],
    "categories": "53 - 59 inch TVs,Televisions,TV & Home Theatre",
    "merchant_id": "210"
  },
  {
    "name_v": [
      "Samsung 55\" The Frame 4K UHD Smart LED TV with Art Mode"
    ],
    "categories": "TV's (55-64\")",
    "merchant_id": "299"
  },
  {
    "name_v": [
      "Samsung 65\" Wood Frame for The Frame TV - Beige Wood"
    ],
    "categories": "Home Theatre Accessories,TV & Home Theatre,TV Frames",
    "merchant_id": "210"
  },
  {
    "name_v": [
      "Samsung 65\" Customizable Frame for \"The Frame\" TV - Oak"
    ],
    "categories": "Other TV & Video Accessories",
    "merchant_id": "299"
  }
]
```

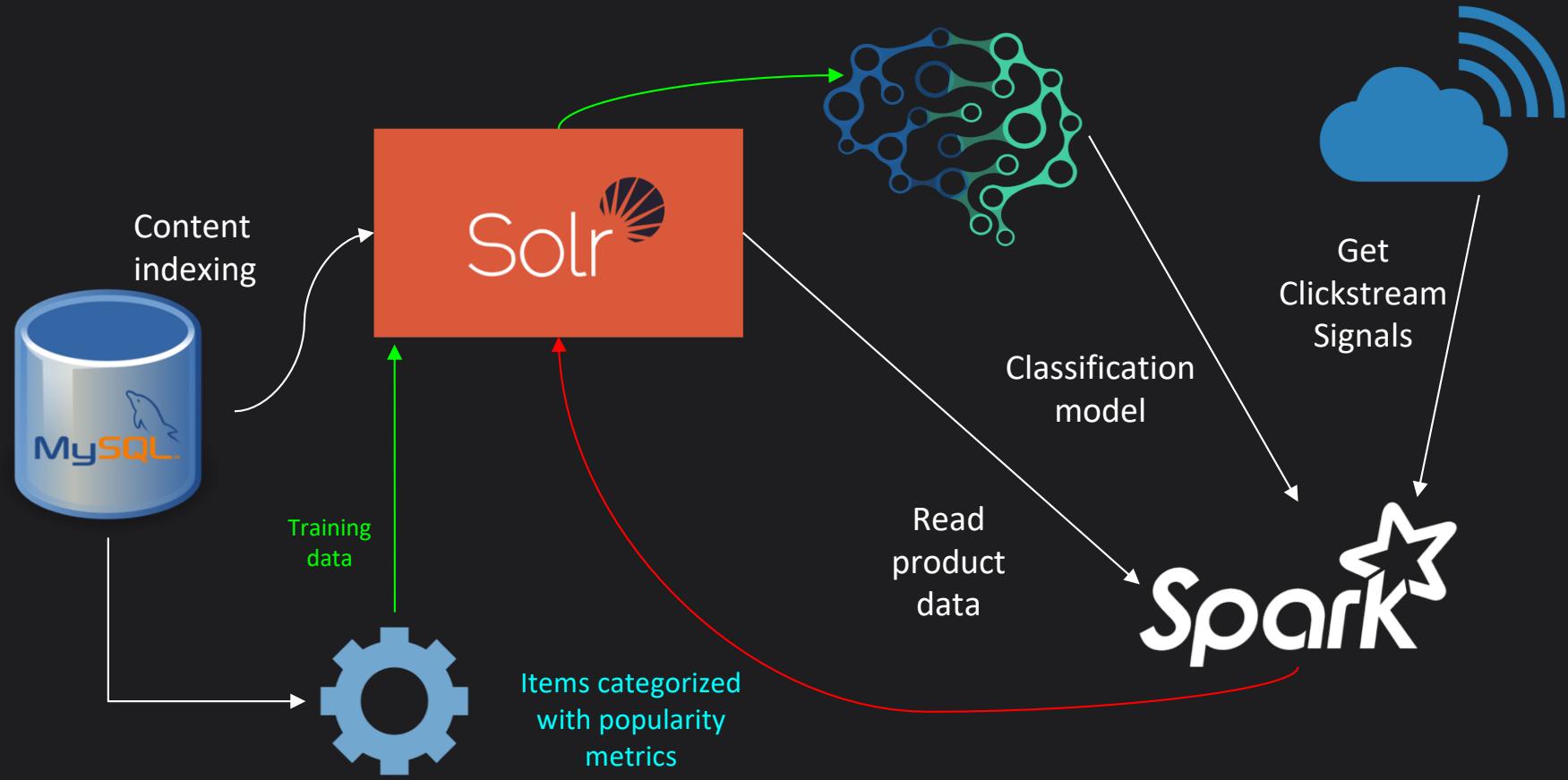
Same products but different category structure

TV accessories that would show up as noise

# Categorization Algorithm



# Content Categorization Workflow



# Voilà!

```
  "docs": [
    {
      "name_v": [
        "Samsung The Frame 55\" 4K UHD HDR LED Tizen Smart TV (UN55LS003AFXZC) - Charcoal Black"
      ],
      "categories": "53 - 59 inch TVs, Televisions, TV & Home Theatre",
      "merchant_id": "210",
      "L1": "Electronics",
      "L2": "Video",
      "item_weight": 1.5928669693999906
    },
    {
      "name_v": [
        "Samsung 55\" The Frame 4K UHD Smart LED TV with Art Mode"
      ],
      "categories": "TV's (55-64\")",
      "merchant_id": "299",
      "L1": "Electronics",
      "L2": "Video",
      "item_weight": 0.6496946375051786
    },
    {
      "name_v": [
        "Samsung 65\" Customizable Frame for \"The Frame\" TV - Oak"
      ],
      "categories": "Other TV & Video Accessories",
      "merchant_id": "299",
      "L1": "Furniture",
      "L2": "Entertainment Centers & TV Stands",
      "item_weight": 4.199292645481367
    },
    {
      "name_v": [
        "Samsung 65\" Wood Frame for The Frame TV - Beige Wood"
      ],
      "categories": "Home Theatre Accessories, TV & Home Theatre, TV Frames",
      "merchant_id": "210",
      "L1": "Furniture",
      "L2": "Entertainment Centers & TV Stands",
      "item_weight": 1.4340761125496102
    }
  ]
```

TVs are now under similar categories with popularity weight

The stands are now correctly classified as 'Furniture' and can be filtered

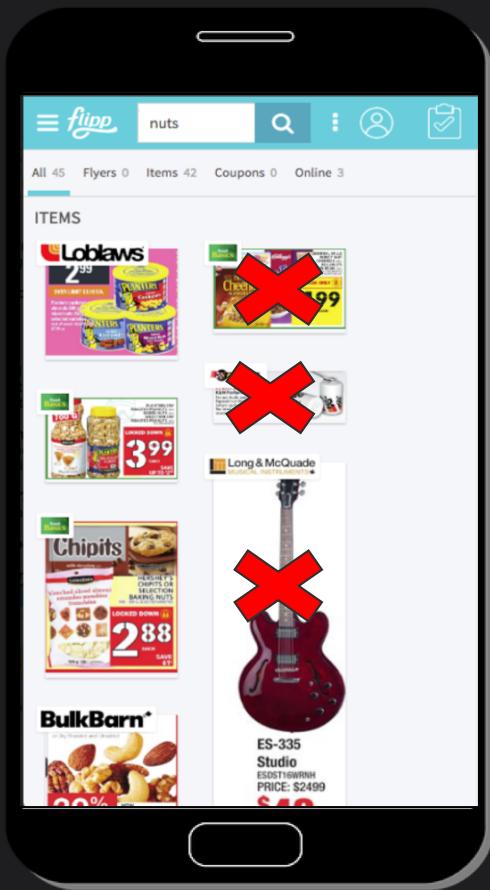
DEEP DIVE

# Understanding User Intent



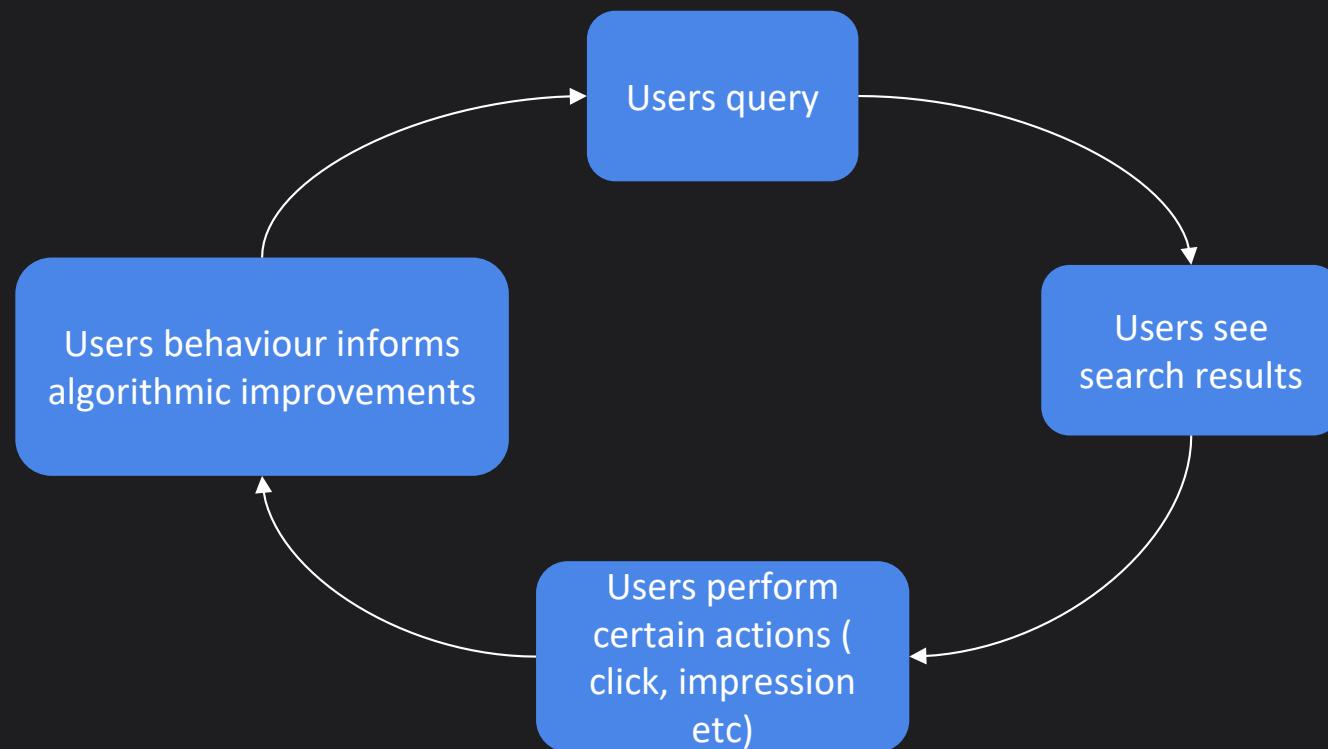
# Why Decipher Query Intent?

A query for “*nuts*”  
ranks *cereal*,  
*guitar* & *car oil*  
*filter* on top due  
to noisy keyword  
signals



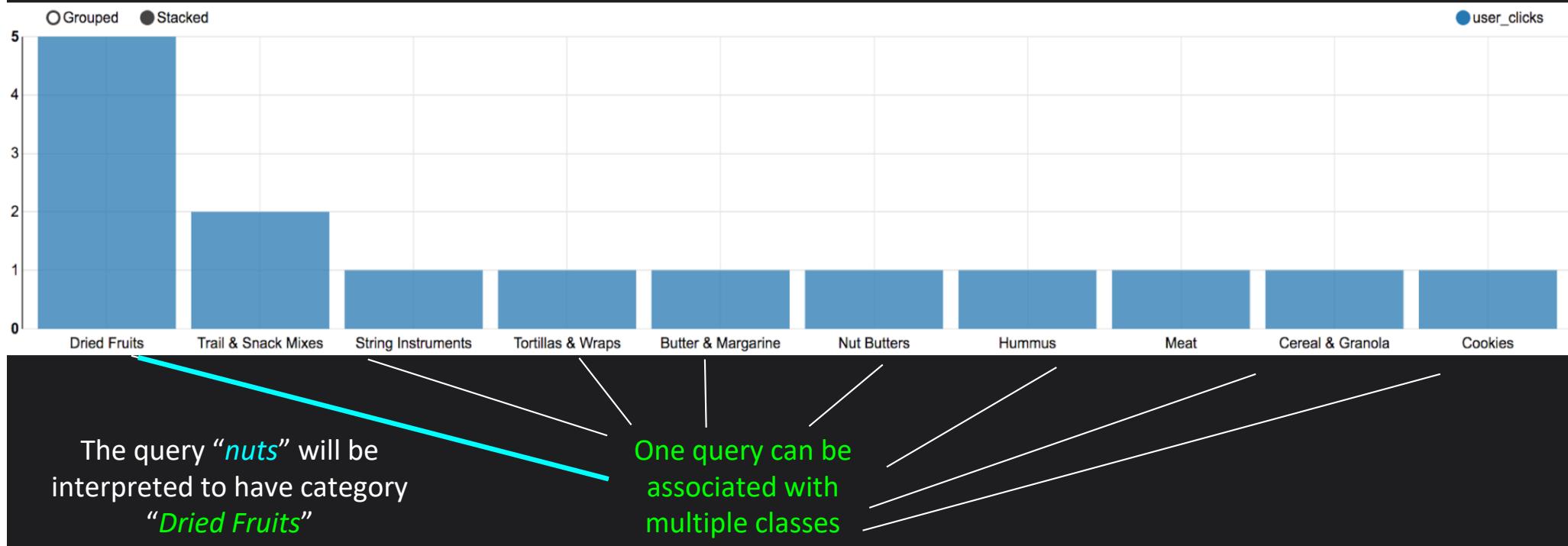
# Reinforced Intelligence

- Using past impression data to improve how latest impression data is interpreted



# Narrowing Down The Query Intent

*What is the probability that a user will click on a particular category class given the query “nuts”?*



*The category with the most clicks for a particular search query would be interpreted as the query category*

# Persisting Query Classifier Model In Solr

Model schema in solr

```
root
|-- category_L1: string (nullable = true)
|-- _L1_conf: double (nullable = true)
|-- category_L2: string (nullable = true)
|-- _L2_conf: double (nullable = true)
|-- category_L3: string (nullable = true)
|-- _L3_conf: double (nullable = true)
|-- category_L4: string (nullable = true)
|-- _L4_conf: double (nullable = true)
|-- category_L5: string (nullable = true)
|-- _L5_conf: double (nullable = true)
|-- category_L6: string (nullable = true)
|-- _L6_conf: double (nullable = true)
|-- category_L7: string (nullable = true)
|-- _L7_conf: double (nullable = true)
|-- isBrand: long (nullable = true)
|-- searchable_query: string (nullable = false)
|-- isMerchant: long (nullable = true)
```

Category classifications for “nuts” variations

```
1 import requests
2 import pandas as pd
3 nut_queries = ['nut thins', 'nuts', 'peanut butter', 'cashews', 'mixed nuts', 'almonds']
4 cats = []
5 for q in nut_queries:
6     cat_frame = {x:y[0] for x,y in requests.get("{}/query?q={}&wt=json&fl=_L1,_L2,_L3,_L4,_L1_prob".format(solr_url,q)).json()["response"]["docs"][[0]].items()}
7     cat_frame["query"] = q
8     cats.append(cat_frame)
9 display(sqlContext.createDataFrame(pd.DataFrame(cats))\
10 .selectExpr('query','_L4 as category','_L1_prob as pred_conf'))
```

▶ (3) Spark Jobs

query	category	pred_conf
nut thins	Cereal & Granola	0.93500465
nuts	Dried Fruits	0.9542686
peanut butter	Nut Butters	0.97368956
cashews	Dried Fruits	0.9877151
mixed nuts	Dried Fruits	0.9845098
almonds	Dried Fruits	0.96738935

# Query Classification Workflow In Solr

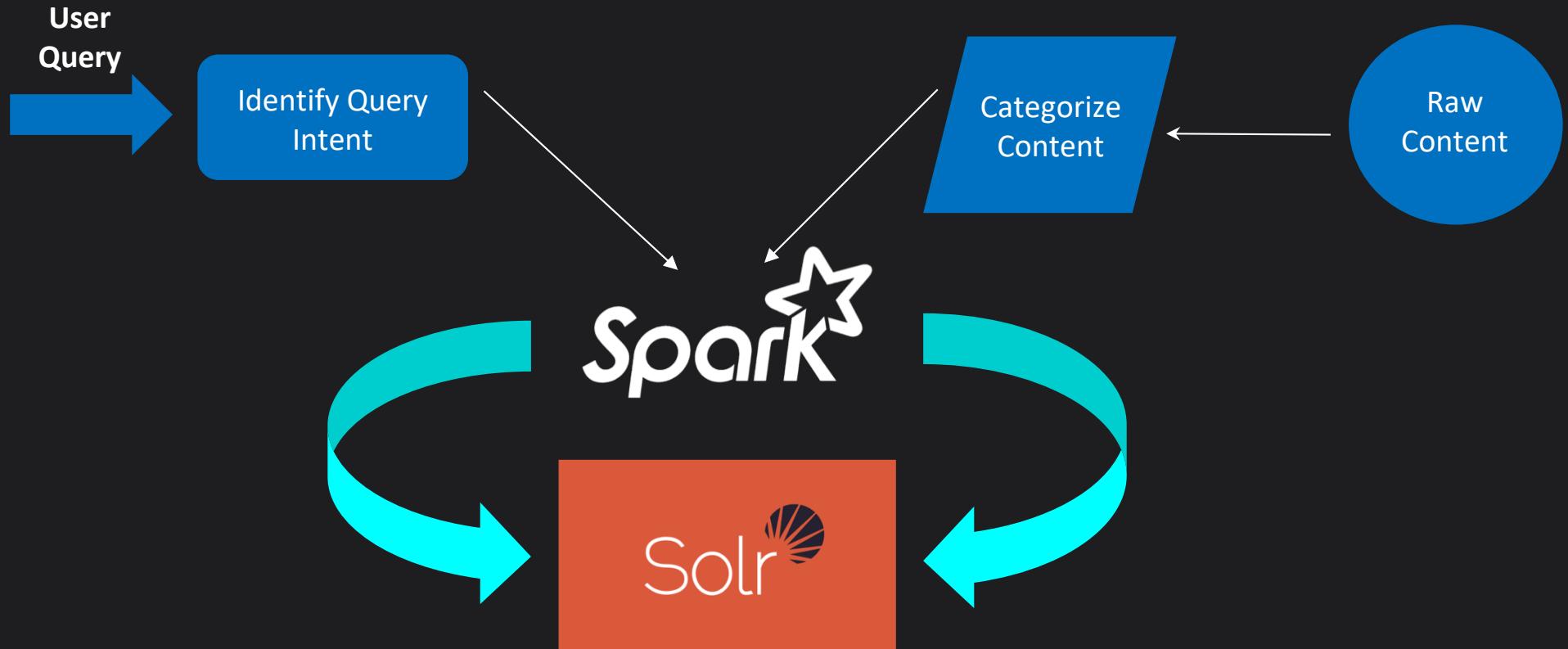
Train

Persist

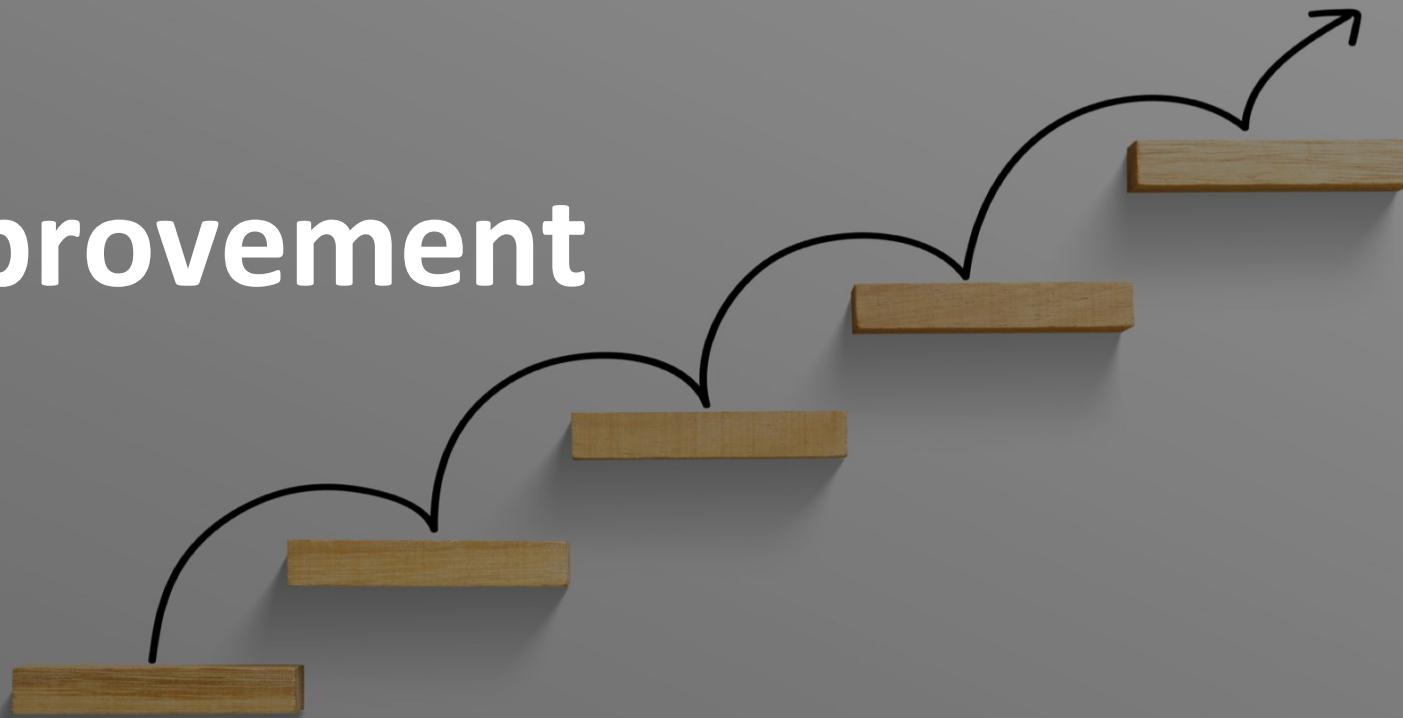
Predict

- Compute category classes
- Classify queries as domain specific entities
- Query as a searchable entity.
- Taxonomic features.
- Query model
- Boost content ranking

# Collaboration!



# Improvement



# Measuring Search Relevance

**Normalized Discounted Cumulative Gain (NDCG)** is a popular method for measuring the quality of a set of search results. It asserts the following:

- **Cumulative Gain:**  $0 \Rightarrow$  Not relevant,  $1 \Rightarrow$  Near relevant,  $2 \Rightarrow$  Relevant

$$\left( CG = \sum_{pos=1}^n relevance_{pos} \right)$$

- **Discounting:** Positional bias

$$\left( DCG = \sum_{pos=1}^n \frac{relevance_{pos}}{\ln(pos+1)} \right)$$

- **Normalization:** Ideal DCG (iDCG)

$$\left( NDCG_{pos} = \frac{DCG_{pos}}{iDCG} \right)$$

# CG Improvement

Q: “Olive Oil”

 <p><b>Cataldi</b> Olive Oil</p> <p>2</p>	 <p><b>Nations Fresh Foods</b> Reginelli Extra Virgin Olive Oil</p> <p>2</p>	 <p><b>None</b> <b>None</b> Olive Oil</p> <p>2</p>	 <p><b>8.99</b> PC® Splendido® extra virgin olive oil 1L</p> <p>2</p>
 <p><b>Marche Adonis</b> RIO MARE Tuna in Olive Oil With Garlic &amp; Red Pepper or Puttanesca Vegetables</p> <p>2/5.00</p> <p>1</p> <p>5</p>	 <p><b>JYSK</b> Olive Oil Waxcloth</p> <p>0</p>	 <p><b>None</b> <b>None</b> Olive Oil</p> <p>2</p>	 <p><b>2.99</b> Reginelli Extra Virgin Olive Oil</p> <p>2</p>
 <p><b>M&amp;M Food Market</b> Rosemary Pull-Apart Bread</p> <p>4.99</p> <p>0</p>	 <p><b>Fortino's</b> CLOVER LEAF TUNA IN OLIVE OIL, 3x80 g OR SOCKEYE SALMON, 213 g</p> <p>3.99</p> <p>1</p>	 <p><b>None</b> <b>None</b> Olive Oil</p> <p>2</p>	 <p><b>5.99</b> Basso Unfiltered Extra Virgin Olive Oil</p> <p>2</p>
 <p><b>Puroliva</b> Extra Virgin Olive Oil</p> <p>\$4.99 ea</p>	 <p><b>7.99</b> Basso Extra Virgin Olive Oil 1L</p> <p>2</p>		

$$2+2+1+0+0+1= 6$$

$$2+2+2+2+2+2= 12$$

# DCG Improvement

Q: “Fish”

	99¢	0.99	Sobeys CLOVER LEAF or OCEAN'S Light Tuna	1.44
	\$2.49 100 G	2.49	Whole Foods Coho Salmon Fillet	1.82
	22.99	22.99	Longos Fresh Wild Halibut Fillet	1.44
	5.99	5.99	Metro DOM RESERVE HOT SMOKED FISH FILLETS	0.62
	8.99 SAVE \$2/LB Hot Smoked Sole Fillet Wildling	8.99	Longos Fresh Wild Caught Sole Fillet	1.11
	\$6.49 REAL DEAL! Reg. \$10.99 400 CALS	6.49	Harvey's Harvey's Meal Deal!*	0

$$1.44+1.82+1.44+0.62+1.11+0= 6.43$$

	2.49	2.49	Coho Salmon Fillet	2.89
	22.99	22.99	Fresh Wild Halibut Fillet	1.82
	5.99	5.99	DOM RESERVE HOT SMOKED FISH FILLETS	0.72
	8.99 SAVE \$2/LB Hot Smoked Sole Fillet Wildling	8.99	Fresh Wild Caught Sole Fillet	1.24
	99¢ Reg. \$10.99 400 CALS	0.99	Clover Leaf or Ocean's Light Tuna	0.56
	\$6.49 REAL DEAL! Reg. \$10.99 400 CALS	6.49	Harvey's Harvey's Meal Deal!	0

$$2.89+1.82+0.72+1.24+0.56= 7.23$$

A photograph of four people in a modern office environment. In the foreground, a man with dark hair and a mustache, wearing a dark blue sweater over a patterned shirt, sits at a desk looking at a laptop. Behind him, three other people are laughing: a woman with blonde hair tied up in a bun, a man in a dark suit, and a woman in a green top. They are all looking towards the laptop screen. The office has large windows showing greenery outside and a blue rug on the floor.

# Spark-Solr: Lessons Learned

# Lessons Learned

- Latency rise in p95 by 40% during indexing



# Lessons Learned

- Read Parallelism
  - Match the number of spark executor nodes to the number of shards in your solr collection.
  - Prefer using solr filter params over `where` clause in Spark.

```
1 val zkhost = "localhost:9983"
2 val signalsOpts = Map("zkhost" -> zkhost, "collection" -> "ecommerce_signals", "query" -> "*:*", "fields" ->
  "count_i, doc_id_s, query_s", "solr.params" -> "fq=userId:[10 TO 1000]", "splits" -> "true")
3 val signals = spark.read.format("solr").options(signalsOpts).load
```

- Indexing improvements
  - Multivalued fields may require extra processing.
  - Avoid using `soft_commit_secs` param in production.
  - Consider using smaller `batch_size`.

# Key Takeaways

- The Spark-Solr framework unlocks a wide variety of machine learning techniques to be applied to search applications.
- Enables faster data analytics in SparkSQL than traditional SQL frameworks.
- Spark-Solr can be used to optimize solr indexing with signal data.
- Storing text classification models in solr has significant semantic benefits over other persistence layers.

# Thank You!