

Deep Learning for Semantic Search in E-commerce



Somnath Banerjee



Head of Search Algorithms at Walmart Labs

<https://www.linkedin.com/in/somnath-banerjee/>

March 19, 2019

Walmart E-commerce search problem

Store Associate



E-commerce Search

provides the functionality of a human but at scale



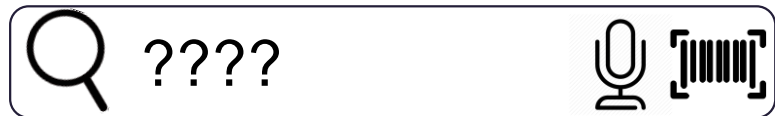
100M+
Customers



100M+
Queries



100M+
Items



Flash Drive
USB Drive
Thumb Drive
Jump Drive
Pen Drive
Zip Drive
Memory Stick
USB Stick
USB Flash Drive
USB Memory
USB Storage Device

Misspelled Queries

Flush Drive
USC Drive
Thamb Drive
Jmp Drive
Pin Drive
Zap Drive
Memory *Steak*
USB *Stock*
USB Flash *Drve*

Q saddle



Horse Saddle

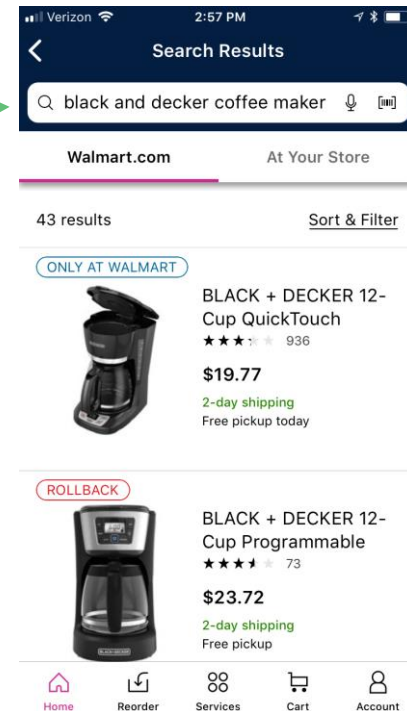
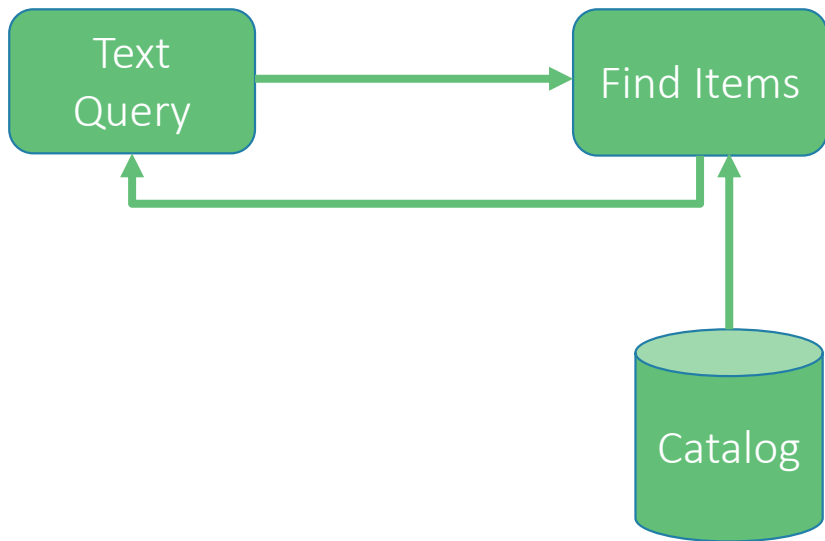


Bike Saddle

Outline

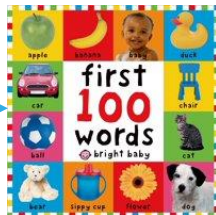
- Core problems of e-commerce search
- Semantic search in e-commerce
- Deep Learning for semantic search
 - Query classification
 - Query token tagging
 - Neural IR
 - Image understanding (sneak peek)

Core of E-commerce Search



Core problems of E-commerce search

Learning book



Tide 100 oz
Tide 100 fl oz
Tide 100 ounce



Levi's
Levi Strauss
Signature by Levi Strauss and Co.



Open vocabulary in query and catalog

Ziploc

product



brand



Ambiguity

Neck style?
Fabric?
No. of pockets?



Missing catalog values

Buying decision is influenced by item attractiveness



Image quality

Presence of expensive items

Tags

Product Name	Price	Original Price	Discount	Rating
Aerosoles Silver Star Pump (Women's)	\$59.99	\$88.99	up to 33% off	4.5 (2)
Easy Street Nancy Pump (Women's)	\$54.95			(0)
Pleaser Amuse 20 (Women's)	\$48.95 - \$65.95	\$48.95 - \$78.95	up to 24% off	4.2 (5)
Stuart Weitzman Marymid Suede Pump (Women's)	\$188.00 - \$281.25	\$375.00	up to 50% off	2.0 (1)
Bella Vita Nara Pump (Women's)	\$72.95 - \$99.95	\$79.95 - \$99.95	up to 9% off	(0)
Easy Street Chiffon Pump (Women's)	\$43.95 - \$49.99	\$49.99 - \$59.99	up to 17% off	4.3 (73)
Walking Cradles Sophia Pump (Women's)	\$49.99 - \$93.95	\$95.00 - \$145.00	up to 50% off	4.3 (4)
Sam Edelman Margie Pump (Women's)	\$119.95 - \$139.95			(0)
Pleaser Seduce 420 (Women's)	\$34.95 - \$49.95			4.5 (33)
Pleaser Pink Label Fab 422 Open-Toe Pump (Women's)	\$57.95			3.5 (2)
Nina Forbes2 D'Orsay Pump (Women's)	\$63.99 - \$84.95	\$85.00	up to 25% off	(0)

Core technical problems of e-commerce search

Matching
query to items

Ranking
items

Pump shoes

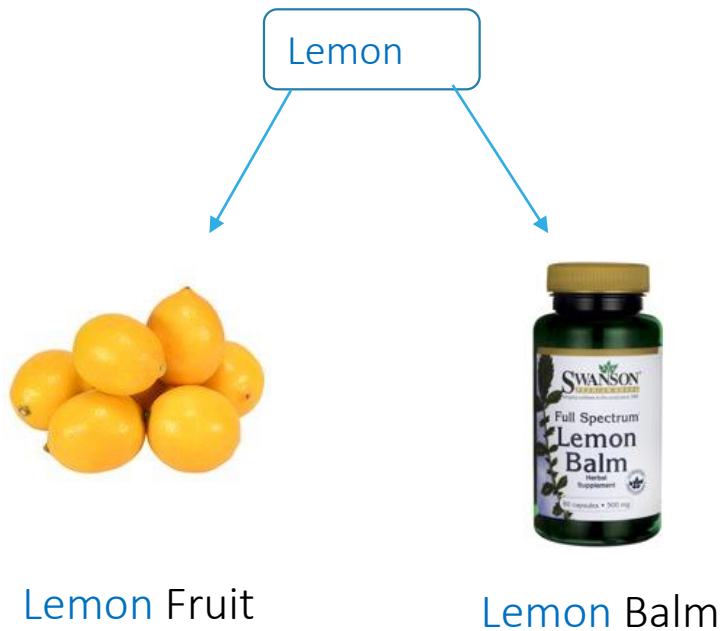


Position 1

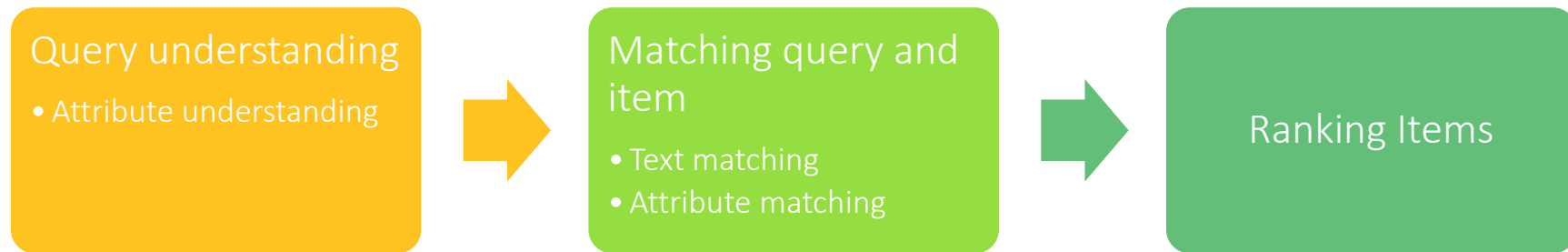


Position 2

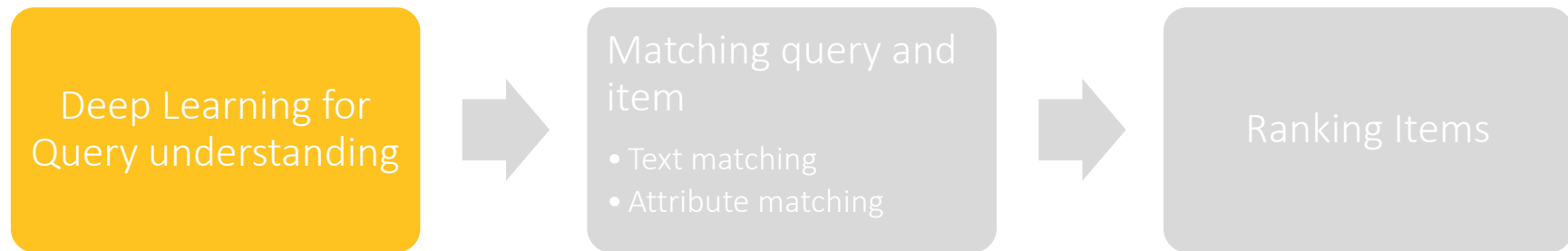
Text matching is not enough



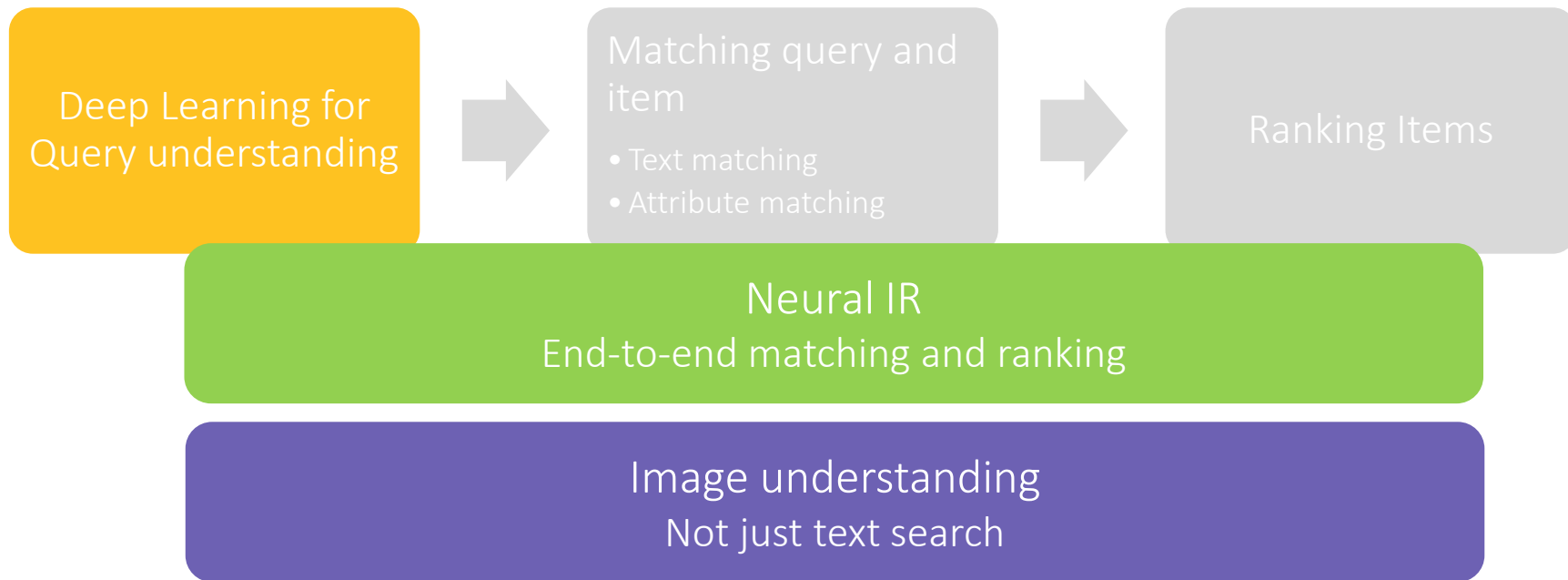
Sematic Search



Deep learning for semantic search



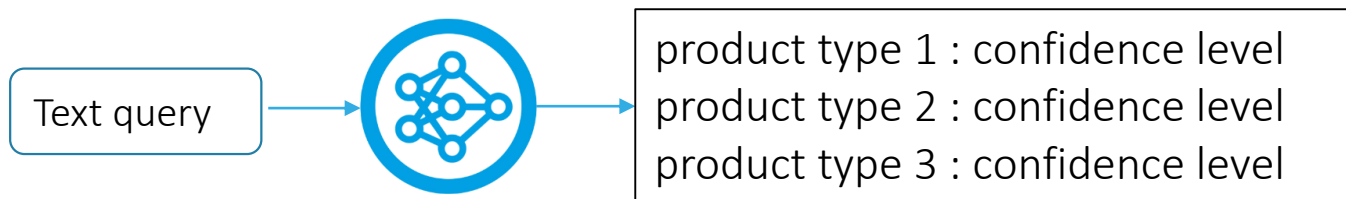
Deep learning for semantic search



Outline

- Core problems of e-commerce search
- Semantic search in e-commerce
- Deep Learning for semantic search
 - Query classification
 - Query token tagging
 - Neural IR
 - Image understanding (sneak peek)

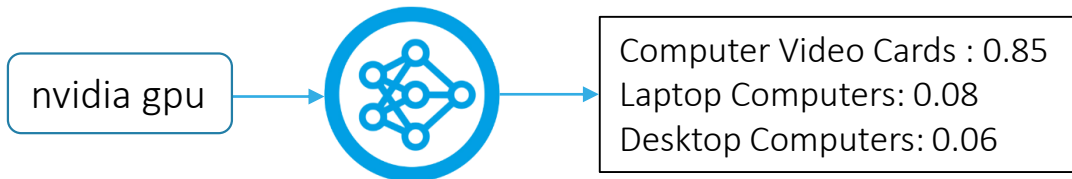
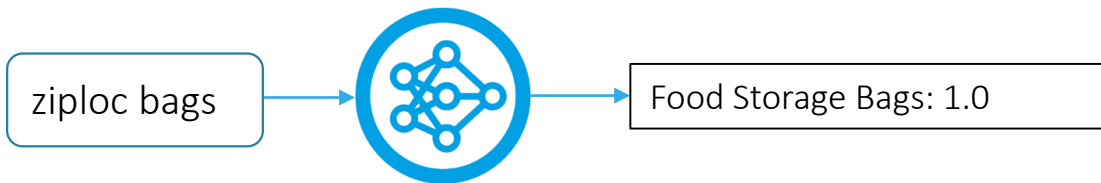
Query Classification



Product Type

- A predefined list
- Indicates a specific product in the catalog
- Every item in the catalog is tagged with a product type

Query classification examples



Hard to balance
precision vs recall

Query classification challenges

Short text

- Queries are of 2-3 tokens

Large scale classification

- Thousands of product types (classes)

Multi-class, multi-label problem

- Same query can have multiple product types

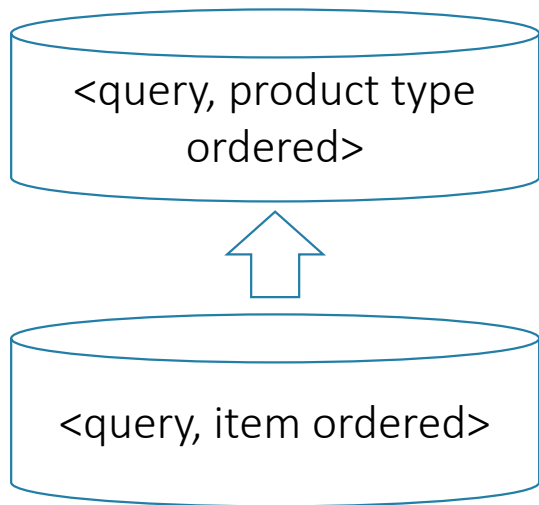
Needs to respond in few milliseconds

- Classifies queries at runtime

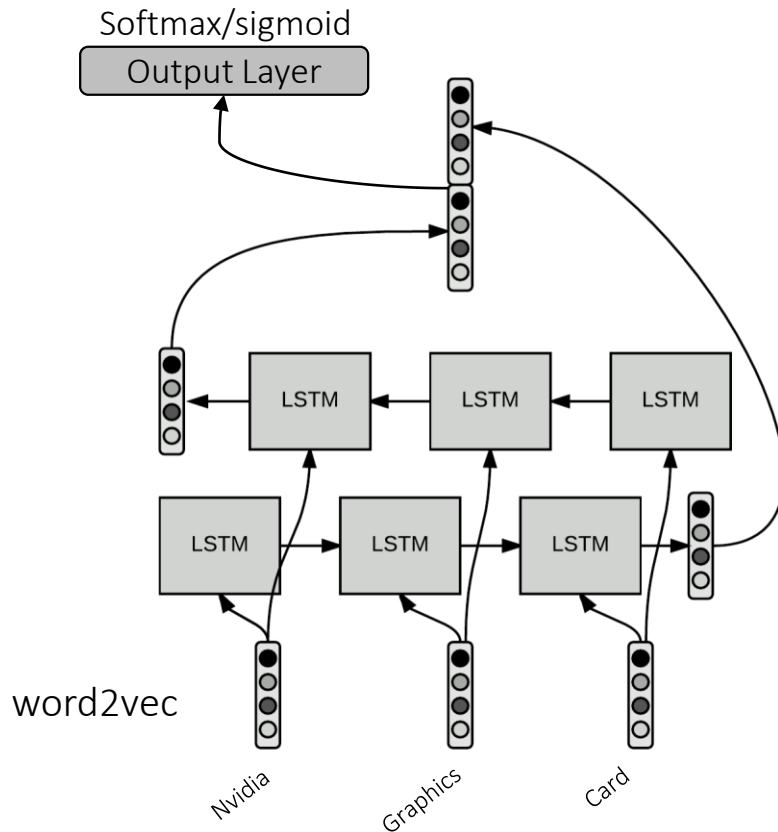
Unbalanced class distribution

- Some product types are much more popular

Data and Model

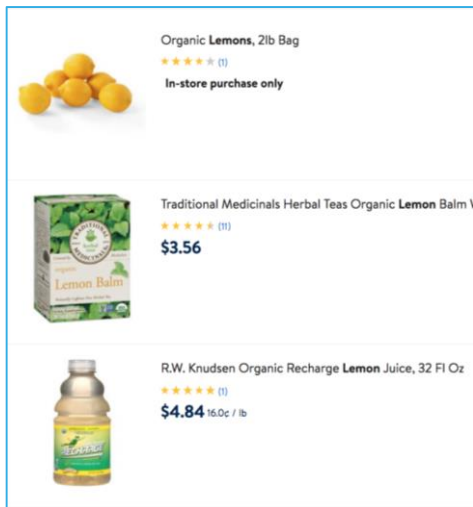


Historical Search Log

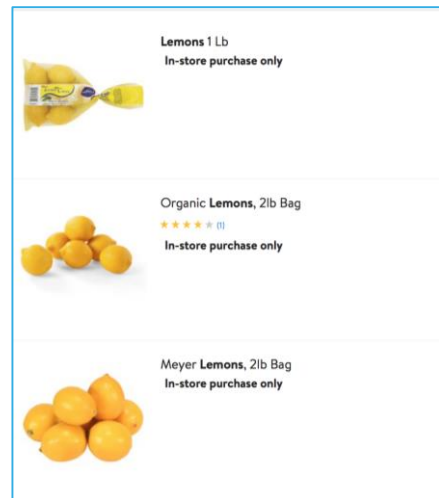


BiLSTM

Usage of query classification



Without Query Classification



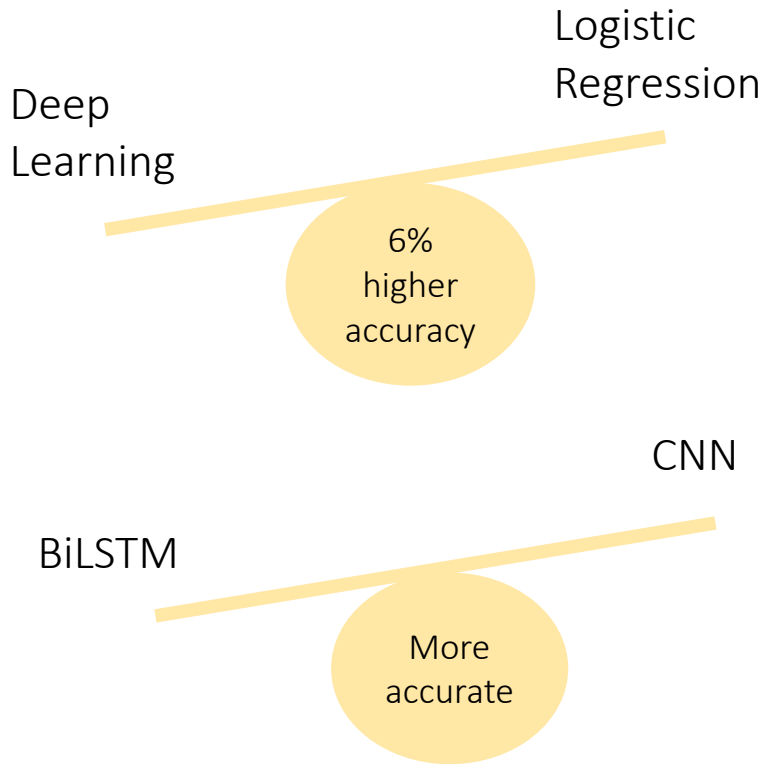
After we understand the query "lemon" as a fruit



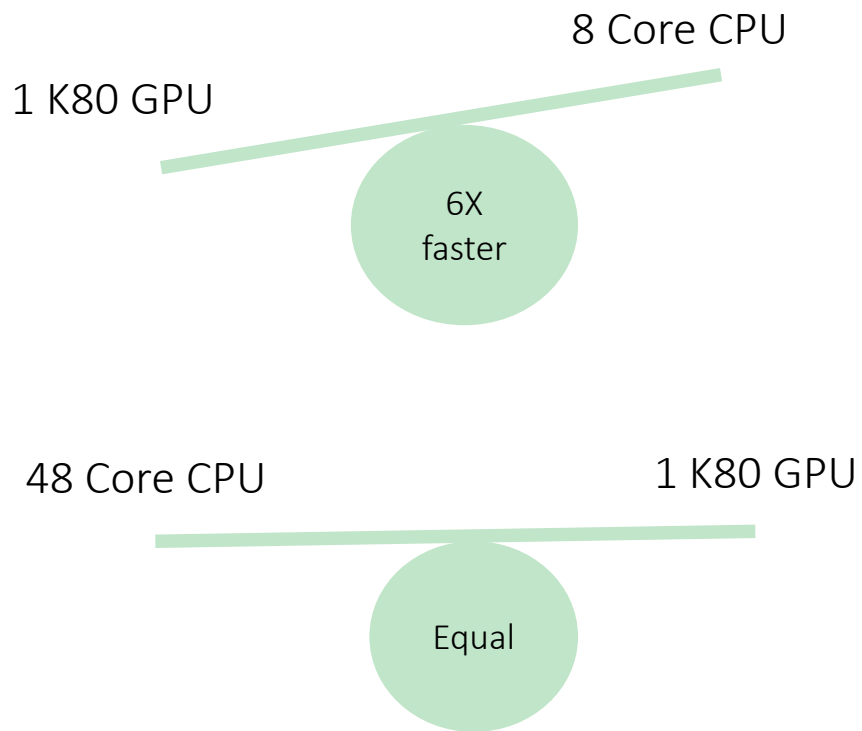
20% reduction of irrelevant items in certain query segments

Key Learnings

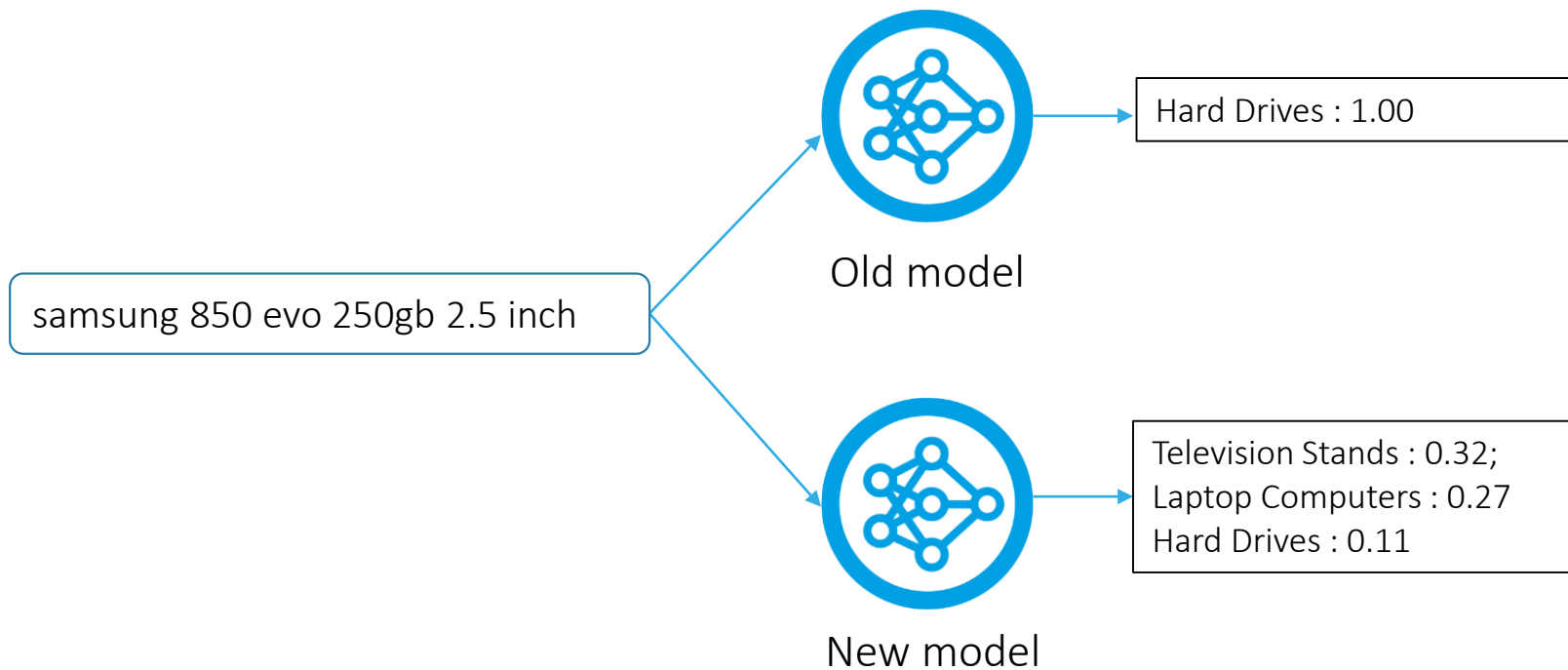
Accuracy



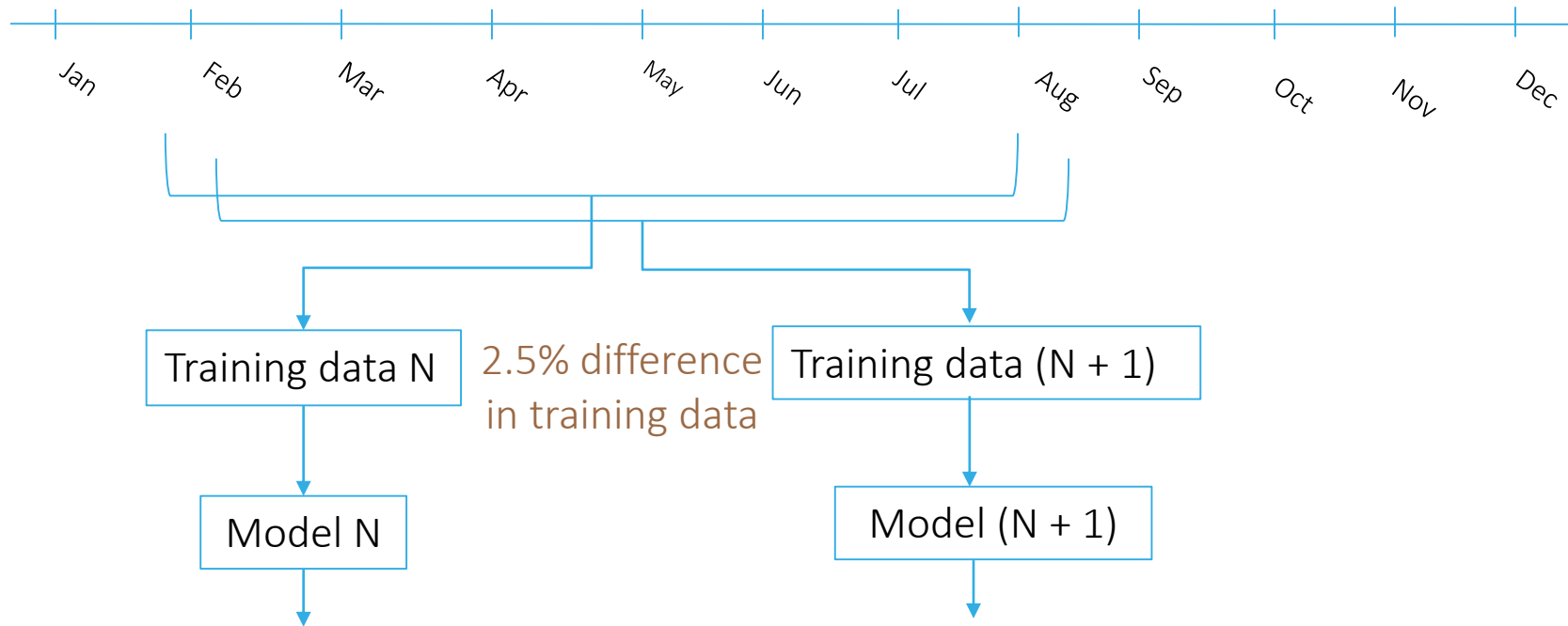
Training Time



Key Learnings - instability in Prediction

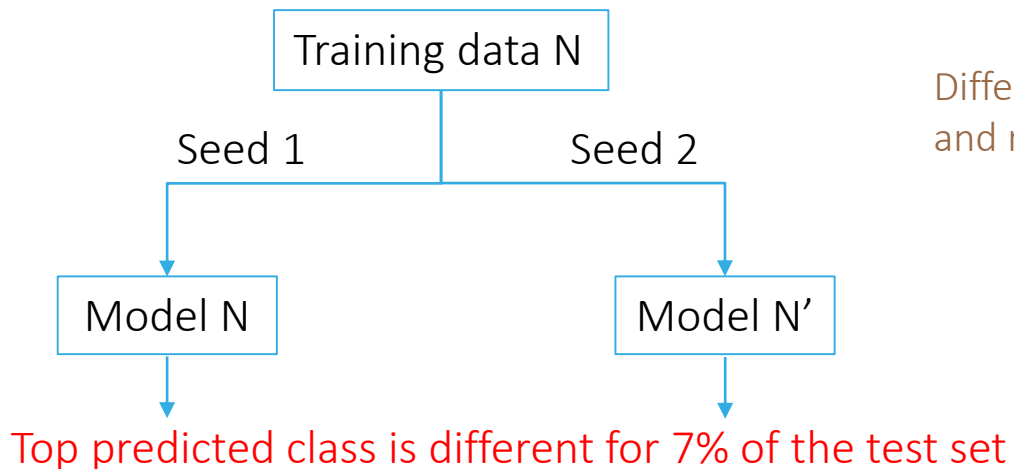


Instability in Prediction



Top predicted class is different for 10% of the test set

Instability in prediction – different seeds



Different tensorflow
and numpy seeds

Sources of Instability

Overfitting

- Deep Learning model has high variance, particularly on the low traffic queries
- Simpler models could be more stable but less accurate

Sigmoid (1-vs-all) classifier is more unstable

- Softmax scores are interdependent across classes and less stable

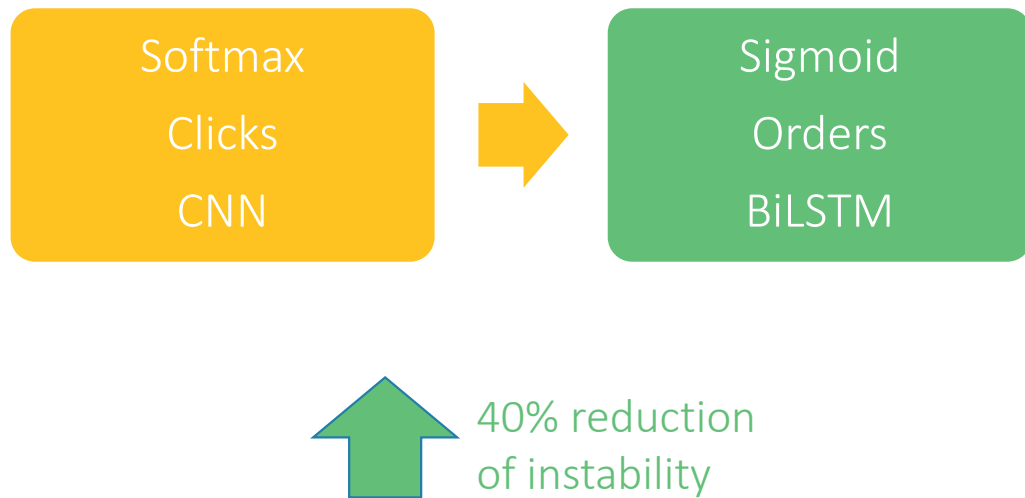
Noisy training data

- Item order data is less noisy than click

Rounding errors in the arithmetic operations

- CPU is more stable than GPU

Reduction of Instability



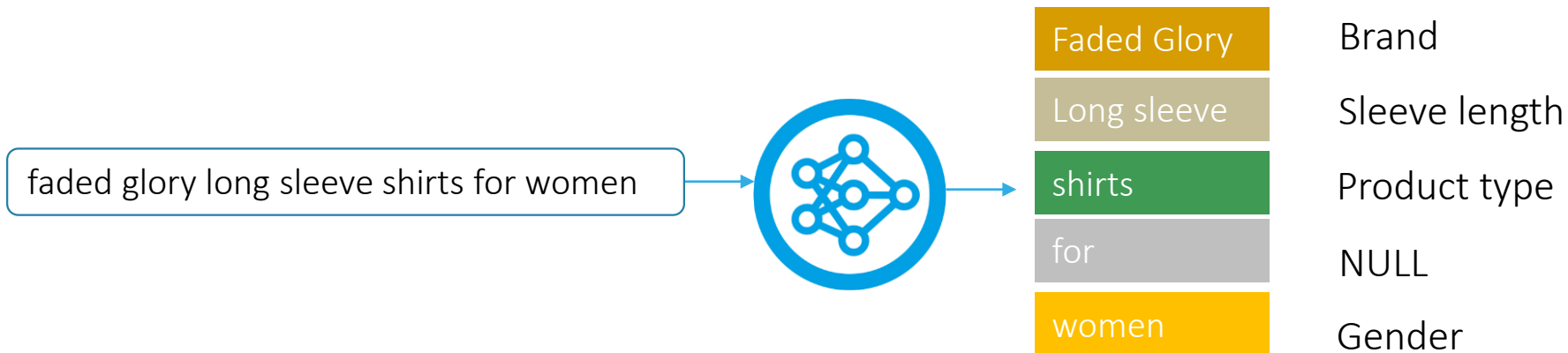
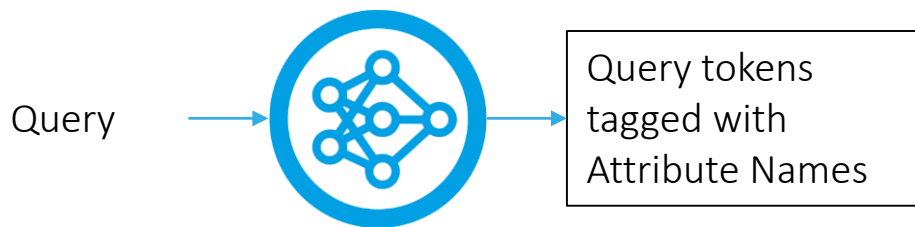
Attributes to match

- Product Type
- Brand
- Color
- Gender
- Age Group
- Size (value & unit)
 - Pack Size
 - Screen Size
 - Shoe Size
 - ...
- Character
- Style
- Material
- ...

Not Feasible – Separate classifier for each attribute

- Too many classes (e.g. 100K+ brand values)
- Sparse attributes; most attribute prediction should be NA
- Creating training data of <query, attribute> is more noisy and inaccurate

Query token tagging



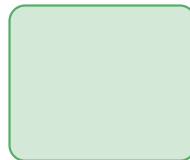
Training data

blue women levis jeans

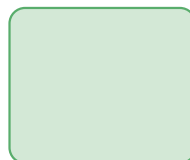
Color Gender Brand Product Type

toys for girls 3 – 6 years

Product Type Gender Age Value Age Unit

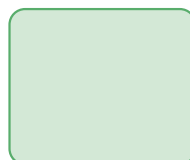


Human curated data

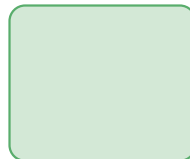


It is a hard task for human

- Is "outside" a product type token in the query, "canopy tents for outside"?

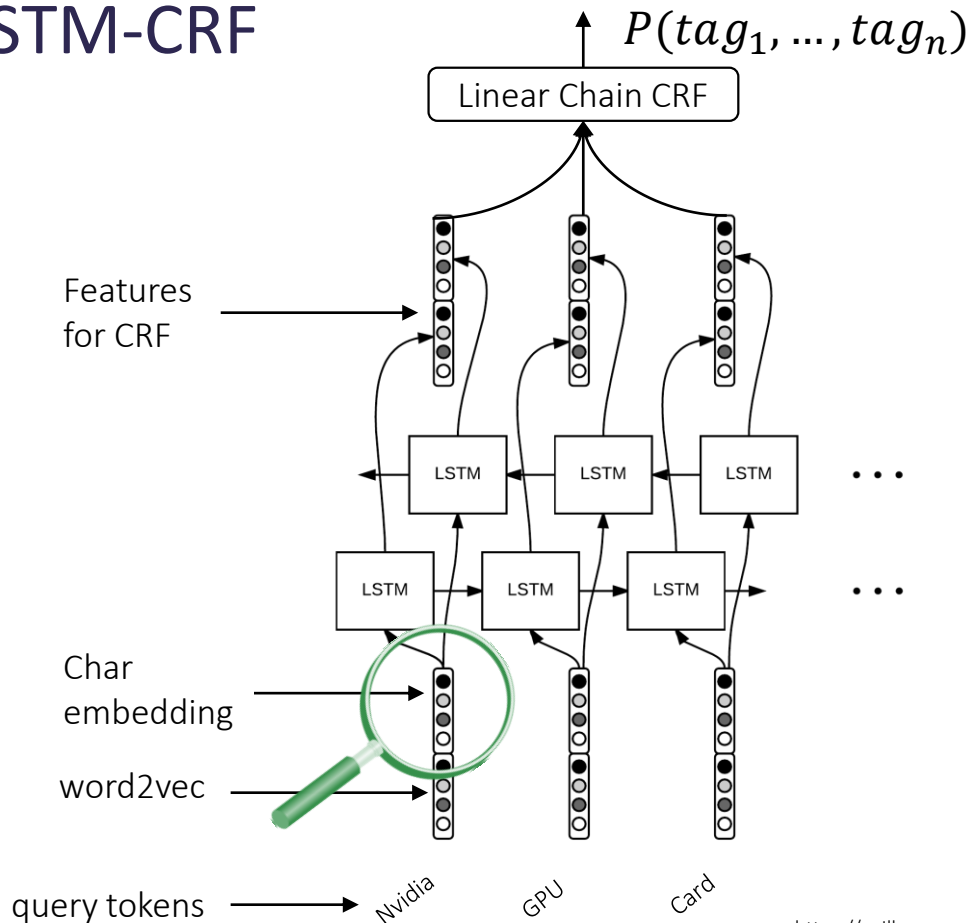


Disagreement between taggers are high (~30%)



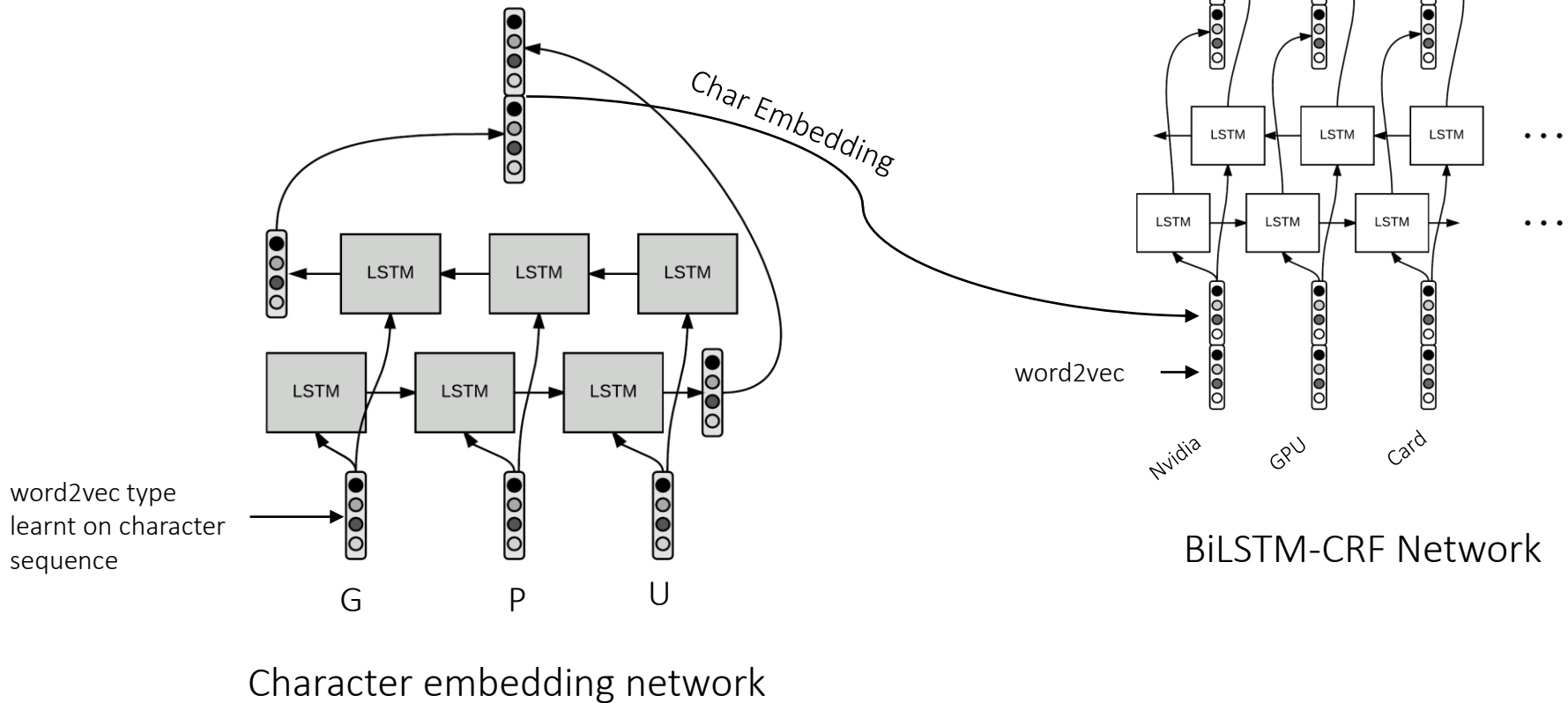
Fortunately 10K training data is a good start

Model – BiLSTM-CRF



<https://guillaumegenthal.github.io/sequence-tagging-with-tensorflow.html>

Char Embeddings



Char Embedding

- Maps a sequence of characters to a fixed size vector
- Handles out of vocabulary words
- Handles misspellings



NULL Product
 Type

Without Char Embedding







Brand Product
 Type





With Char Embedding

Improving search results using query tagging

Q Women citizen eco drive watch  

	<p>Citizen Eco-Drive Titanium Perpetual Atomic Mens Watch AT4010-50E</p>	<p>\$398.27 #1</p>
	<p>Citizen Eco-Drive Blue Angels Chronograph Atomic Mens Watch, AT8020-03L</p>	<p>\$395.95 #2</p>
	<p>Citizen Eco-Drive Promaster Diver Stainless Steel Mens Watch BN0191-55L</p>	<p>\$218.04 #3 from Watchsavings</p>
	<p>Citizen Eco-Drive Skyhawk Blue Angels A-T Perpetual Mens Watch JY8058-50L</p>	<p>\$427.93 #4</p>

Before

	<p>Citizen Eco-Drive L Sunrise Women's Watch, EM0320-59D</p>	<p>\$490.00 #1 from Top One International Corp.</p>
	<p>Citizen Women's EW1544-53A Gold Stainless-Steel Plated Eco-Drive Fashion Watch</p>	<p>\$117.99 #2 from AreaTrend</p>
		
		

Regex match will be incorrect for queries like

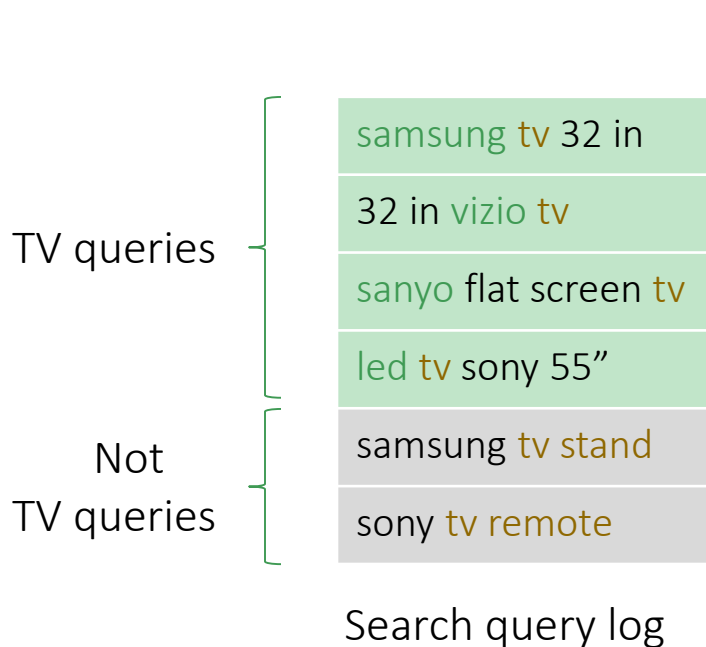
pioneer women dinnerware

wonder women bedding

spider man car seats

After understanding the Gender token

Other use cases of query tagging



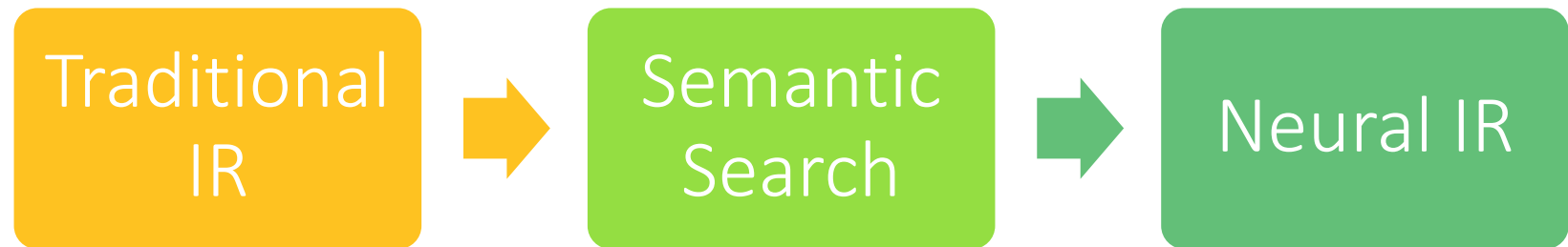
Customer Demand Analysis

- Most searched brand of TV

Attribute filter suggestion

- Suggest top attributes (e.g. brand, screen size) that customers look for for in a product type query (e.g. TV)

Neural IR

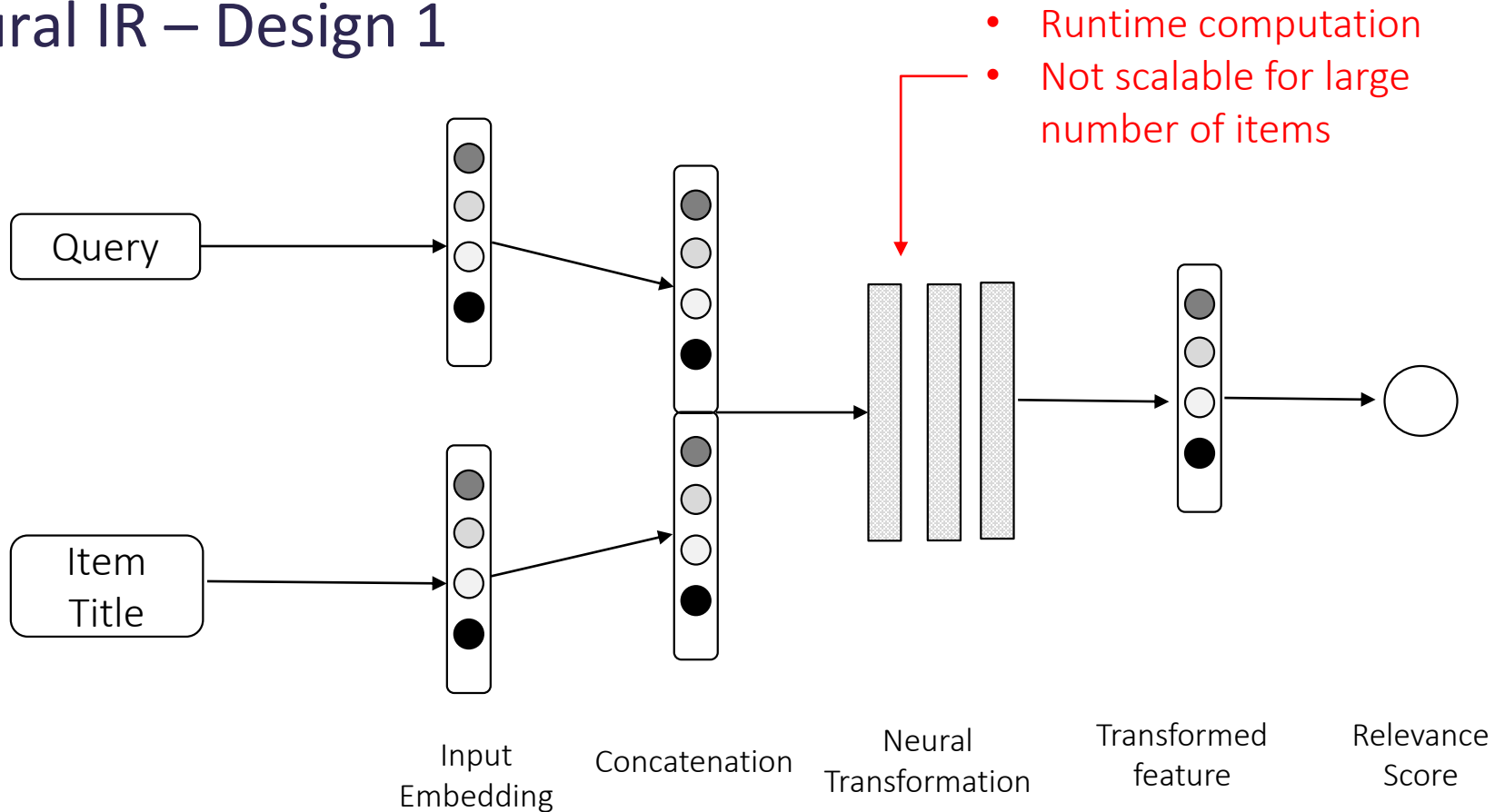


Token and synonym
match
Learning to Rank

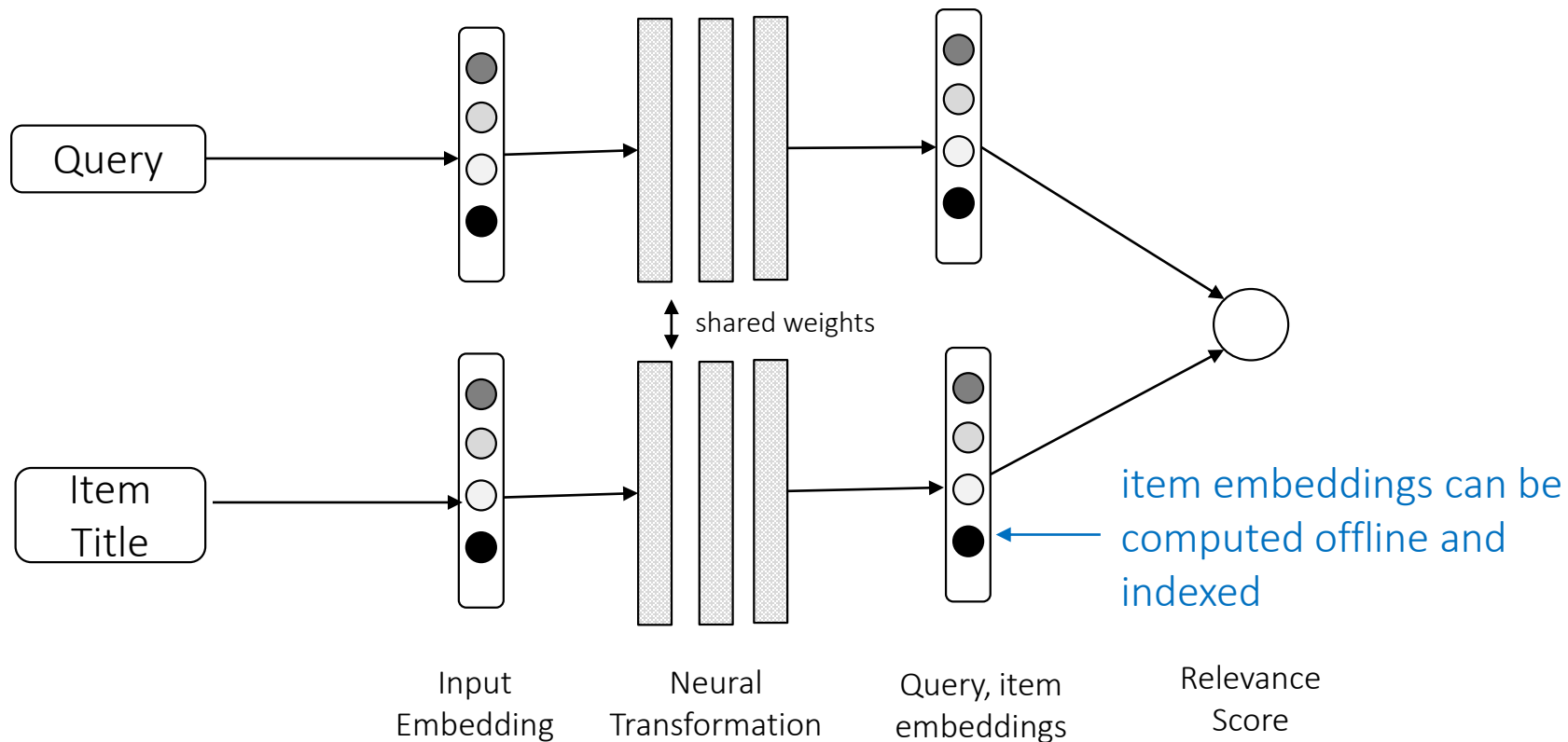
- Attribute extraction
- Token, synonym and attribute match
- Learning to rank

End-to-end matching
and ranking

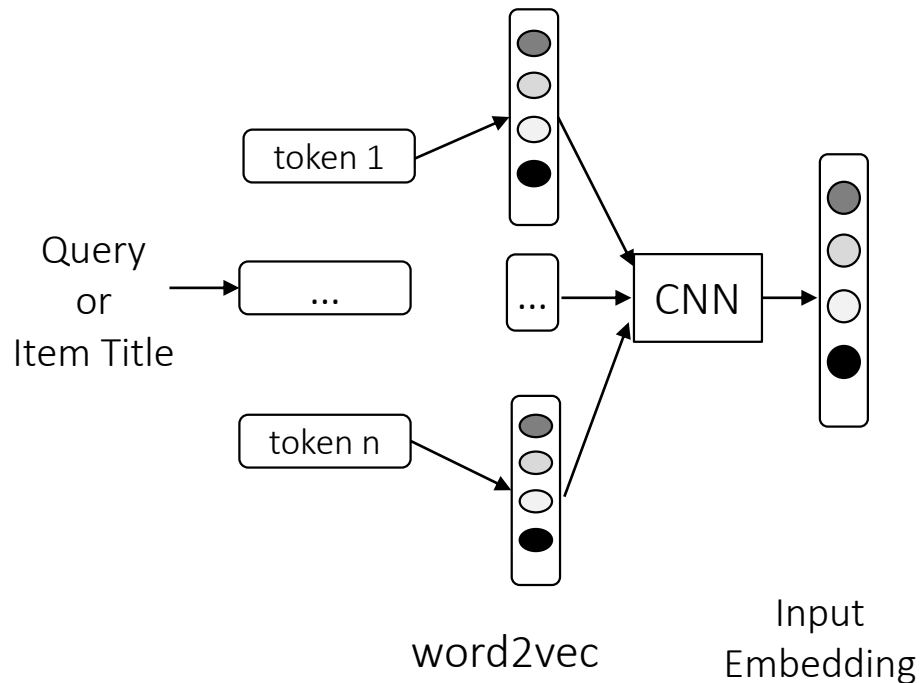
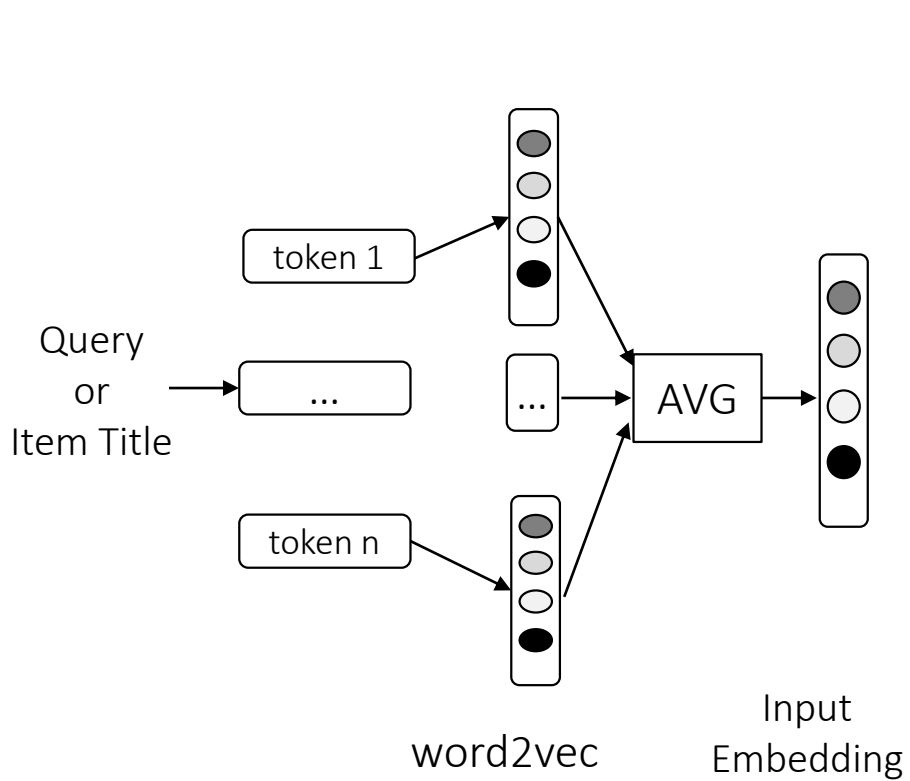
Neural IR – Design 1



Neural IR – Design 2



Input Embedding



Comparable Accuracy

Training Data

query, item title, click through rate (ctr)*

Historical search log

*Position bias correction for ctr of a query, item pair

$$ctr = \frac{\sum_r clicks_corrected_r}{\sum_r impressions_r}$$

$$clicks_corrected_r = clicks_r + (impressions_r - clicks_r) * P(click | r)$$

r = rank at which the item was displayed

Training Loss

Point-wise

$x_q = \text{query features}$

$x_p = \text{item features}$

$f(x_q, x_p) \rightarrow \text{ctr}$

Regression problem

Sigmoid cross entropy loss

Pair-wise

x_q

Brooks shoes

query

x_p



relevant

x_n



less relevant

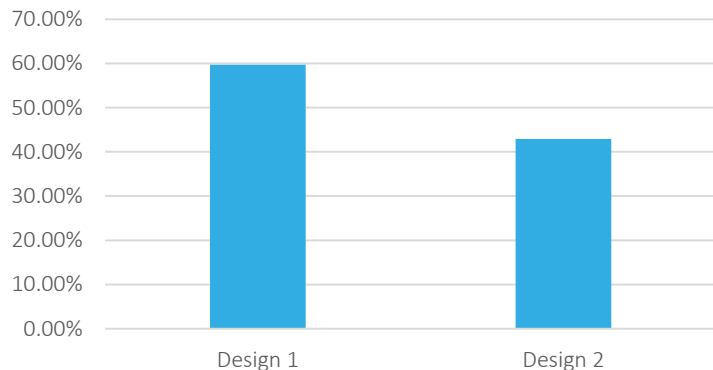
$f(x_q, x_p) > f(x_q, x_n)$
when $\text{ctr}(x_q, x_p) > \text{ctr}(x_q, x_n)$

Minimize pair inversions

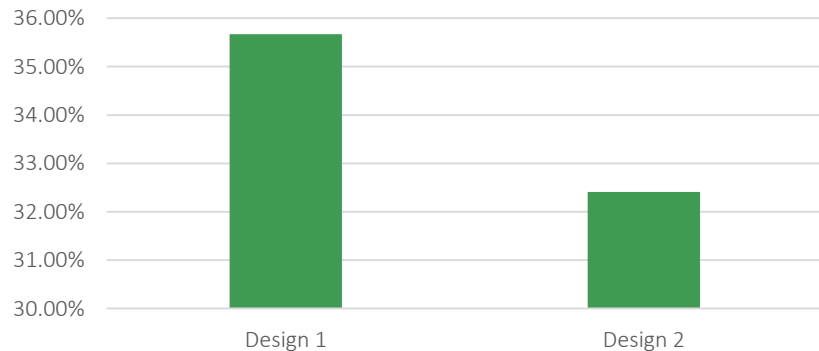
Pair-wise logistic loss

Accuracy on pair-wise loss

NDCG@10 lift against baseline



Pair Accuracy lift against baseline



NDCG captures quality of overall ranking

Pair accuracy captures if higher ctr (relevant) items ranked above the lower ctr items

Neural IR

Pros

- End to end approach
- Enables Semantic matching implicitly
- Handles different data types (text, image)

Cons

- Not scalable (yet)
- Not so successful (yet)

Image understanding

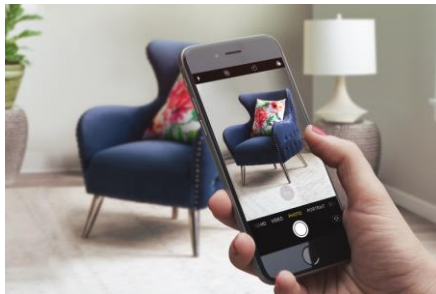
Attribute Prediction



Predicted Attributes

- Product type
- Style
- Material
- Color

Visual Search



Compatible Outfit



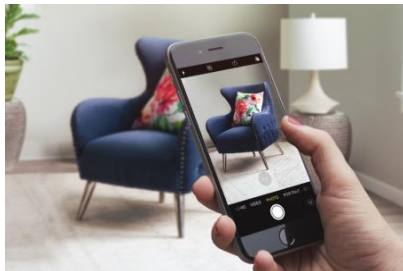
Image understanding key learnings

Attribute Prediction



- Multi-task learning is more accurate
- Predicting style is harder than predicting product type

Visual Search



- A/B test on hayneedle.com
- Comparable results against a well established startup

Compatible Outfit



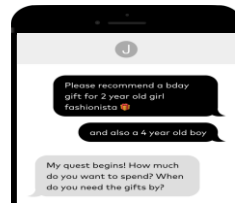
- Under exploration
- Early results beating token based approach

Future



Evolution of mobile phone

Future



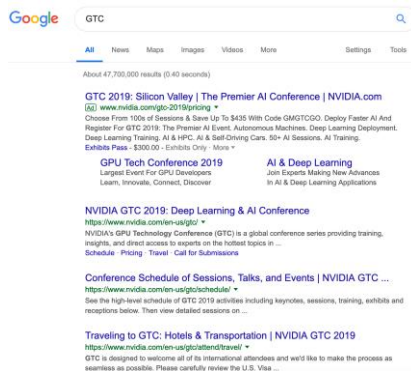
Conversational commerce



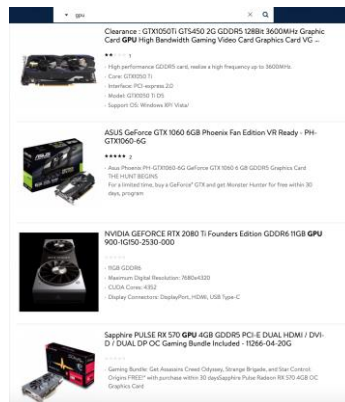
V-Commerce



Seamless search and personalized results



Web Search



E-commerce Search

Thank You