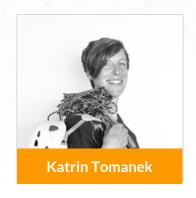


# **Content-Driven Commerce @VigLink**

text-analytics for content monetization at scale







## Introduction



### This talk is about

- Business and Architecture
- Link Insertion in Depth
- Product Categorization

# Affiliate Marketing



```
http://affiliate.acme.com/?
    sub_id_1=<publisherId>&
    sub_id_2=<clickId>&
    url=encode(http://acme.com/products/t8000-headphones)
```

# **LINK** INSERTION



#### Why I Love Photography

APRIL 12, 2014 | 6 COMMENTS | SHARE THIS ARTICLE

The V3 has been completely re-designed. Now the camera has the ability to shoot 120FPS HD video in slo-motion at 1280X720 resolution. The Fujifilm FinePix S4200 Digital Camera is a great camera to try. The one thing they kept with the 1 system and improved upon is indeed the SPEED.



Canon EOS Rebel T3i 18.0 MP DSLR Camera - Black -EF-S 18-55mm IS ...

The <u>Transcend 16GB Class 10 SDHC Flash Memory Card</u> is a great memory card to use for your Fugi. The GOOD news I guess is that this is an all new Nikon 1. They are not dropping the line but instead they beefed it up for even better video capabilities, speed and also packed it with a tilt EVF and a new 18 MP sensor. The You CAN add an external EVF but that always just adds a hump, which these days there is NO reason for. Cameras today are FINALLY getting away from the add-on EVF humps, so why Nikon ditched their internal EVF is a mystery to me. You can now pre-order the V3 with the new 10-30 PD lens for \$1196.95 at 8&H Photousing THIS link

#### CATEGORIES

Technology

Brand

Design

Video

Music

AD

RECENT PHOTOS

### Business



- Founded in late 2008
- Profitable Q4 2013
- \$20 million Series C in 2014
- Emergence, Google Ventures, First Round Capital, RRE,
   Foundry, Silicon Valley Bank
- 45 employees

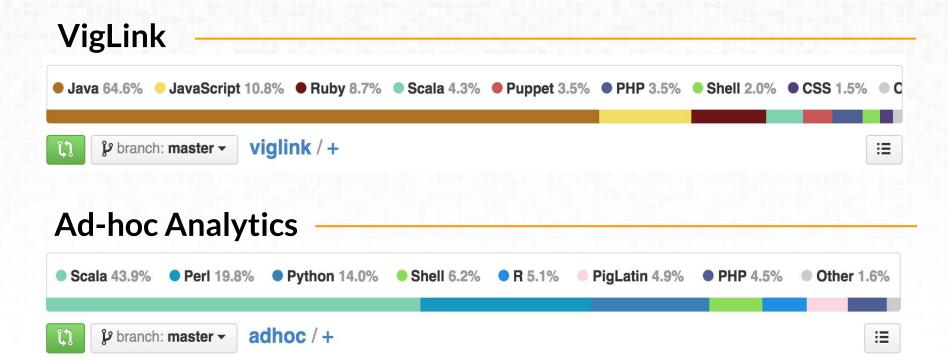
### **Metrics**



- Hundreds of thousands of publisher campaigns
- ~50k merchants (eBay & Amazon + networks like
   Commission Junction)
- > 10 billion page views per month
- Millions of monetized clicks/month
- Growing portion of clicks from Link Insertion

## Programming Language Usage





# High-level AWS Architecture

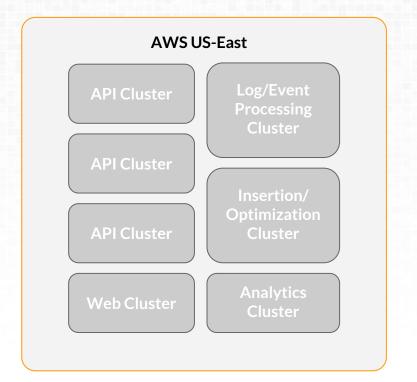


**AWS US-West** 

**API Cluster** 

**API Cluster** 

API Cluste

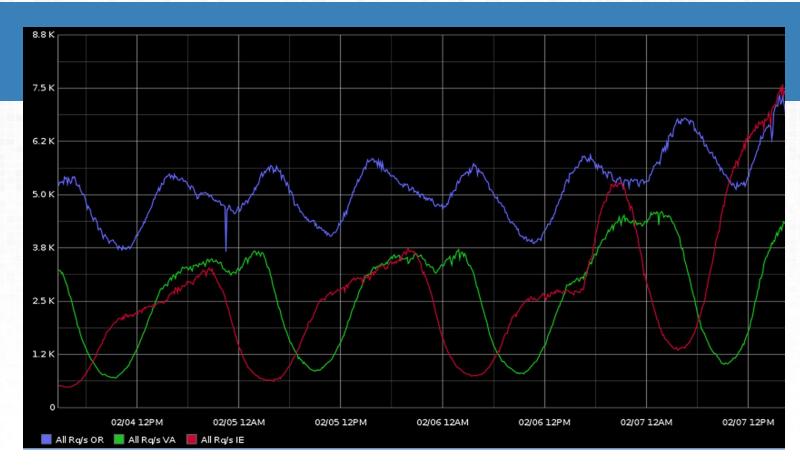


AWS EU-West

API Cluster

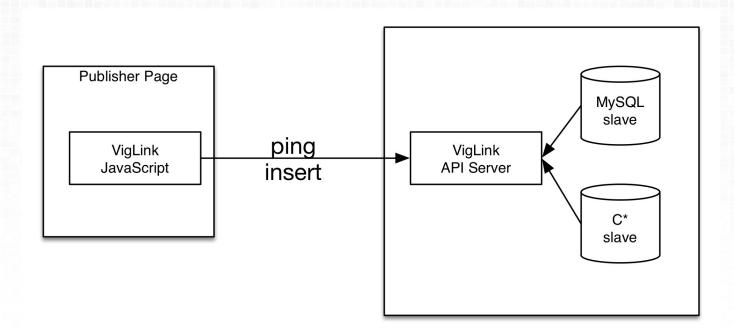
API Cluster

### RPS BY DC



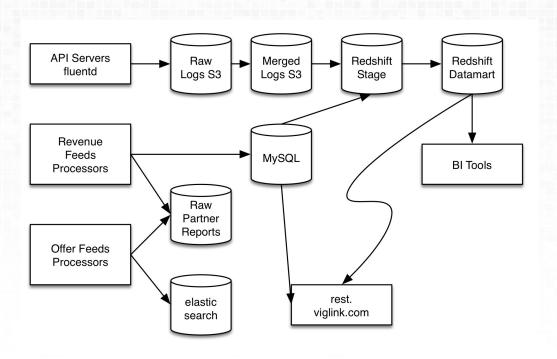
# JavaScript and API





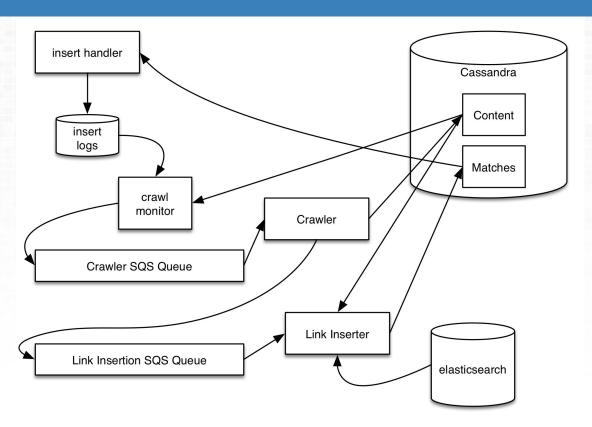
# Log and Feed Processing





# **Crawling and Link Insertion**





# **Link Insertion in Depth**





#### Why I Love Photography

APRIL 12, 2014 | 6 COMMENTS | SHARE THIS ARTICL

The V3 has been completely re-designed. Now the camera has the ability to shoot 120FPS HD video in slo-motion at 1280X720 resolution. The Fujifilm FinePix S4200 Digital Camera is a great camera to try. The one thing they kept with the 1 system and improved upon is indeed the SPEED.



MP DSLR Camera - Black -EF-S 18-55mm IS ...

The <u>Transcend 16GB Class 10 SDHC Flash Memory Card</u> is a great memory card to use for your Fugi. The GOOD news I guess is that this is an all new Nikon 1. They are not dropping the line but instead they beefed it up for even better video capabilities, speed and also packed it with a tilt EVF and a new 18 MP sensor. The You CAN add an external EVF but that always just adds a hump, which these days there is NO reason for. Cameras today are FINALLY getting away from the add-on EVF humps, so why Nikon ditched their internal EVF is a mystery to me. You can now pre-order the V3 with the new 10-30 PD lens for \$1196.95 at 88H Photo using THIS link.

#### CATEGORIES

Technology

Brand

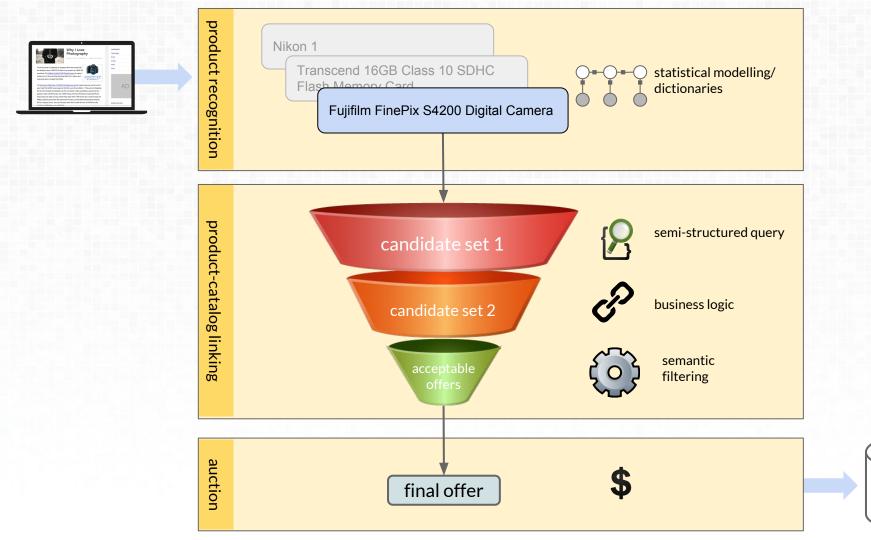
Design

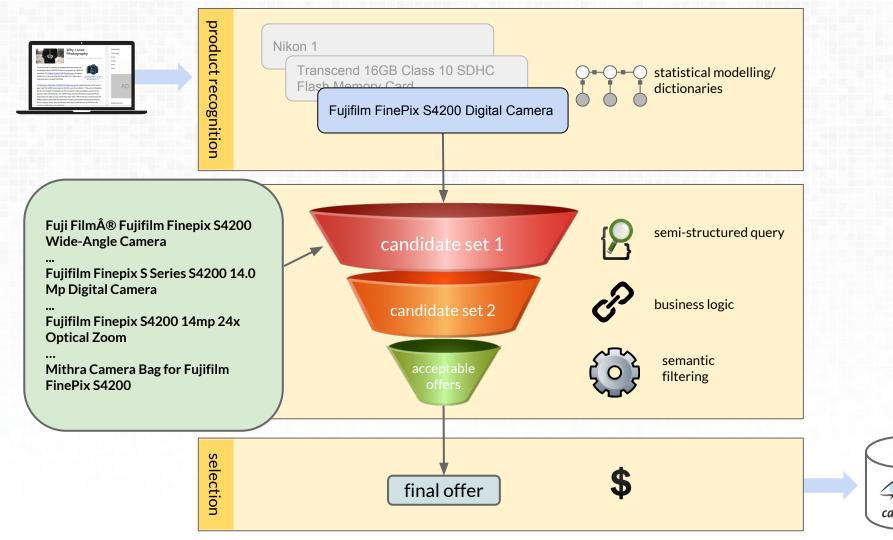
Video

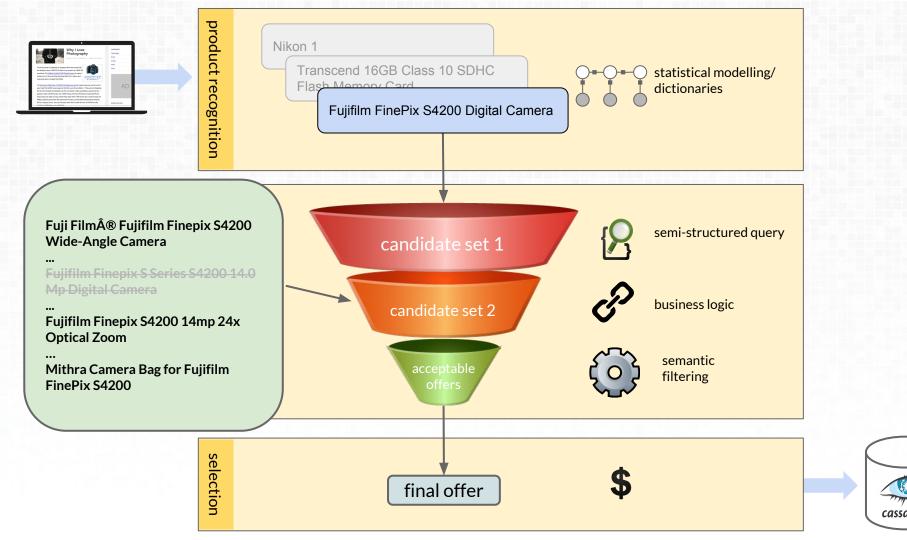
Music

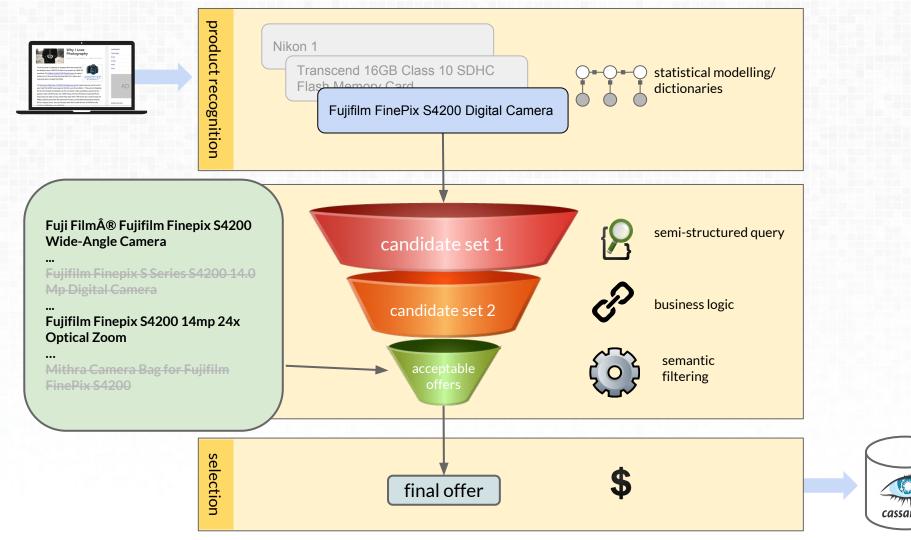
AD

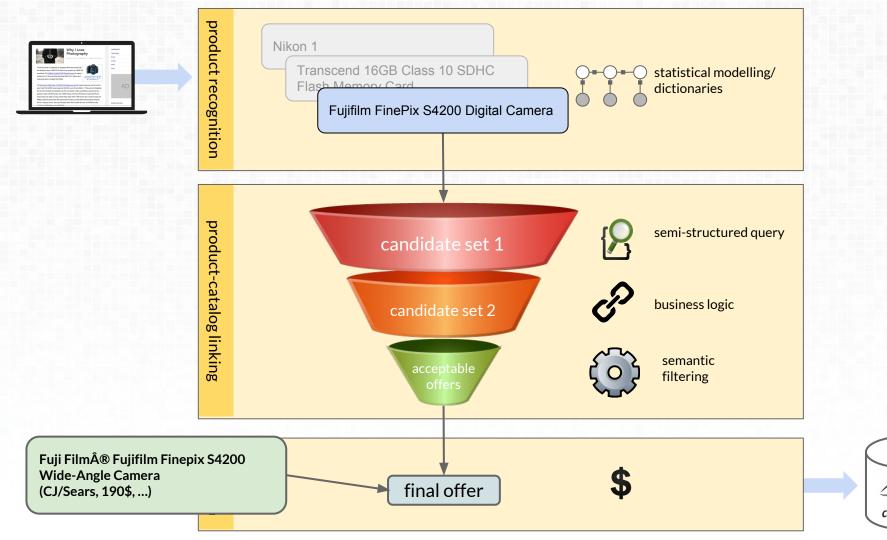
RECENT PHOTOS











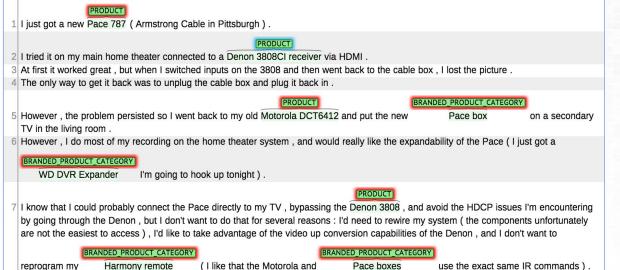
## NLP - related tasks



- product recognition
- linking product to offers
- offer database creation
- understanding offers
- taxonomy

# **Product Recognition**





#### product-like mentions

- products
- product lines
- (branded) product categories

## **Product Recognition**



- typical NER problem
  - sequential tagging
  - impossible to keep all product terms in dictionary
    - spelling variations
    - new products, sub-products...
    - ambiguity, context-dependent
- hybrid approach
  - statistical models
  - dictionary
  - rule-based
  - $\circ \rightarrow ensembles$ 
    - thresholding (R-P lever)

## Sequence Labeling for Product Recognition



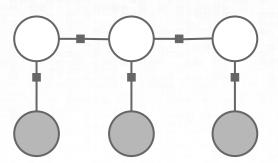
#### Conditional Random Fields

sequence labeling model (linear chain CRF)

	0	0	B-P	I-P	I-P	0	0	0	0	B-BPC	I-BPC	
	my	old	Motorola	DCT	6412	and	put	the	new	Pace	box	

#### features

- word identity
  - token (stemmed)
- orthographical/morphological
  - word pattern
  - numbers, dashes, upper/lower case
- contextual
  - features of neighboring tokens
  - word modifiers (pronouns, articles, quantifiers...)
- lexicon lookups
  - brand names, typical product features, etc



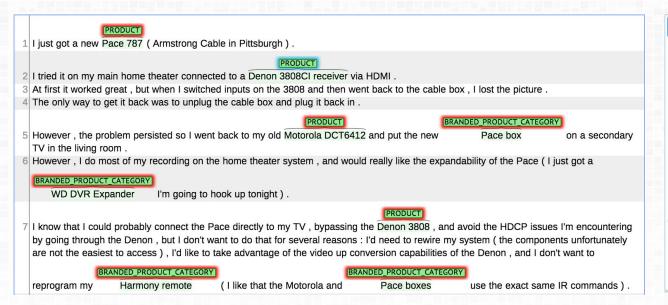
# Training (& Test) Data



- CPROD 1 & 2
  - manually annotated data

	AU	CE	FS	misc
docs	709	2006	156	76
# entities	465	1572	724	1008
% Р	74%	77%	36%	76%
% PL	18%	16%	57%	10%
% BPC	8%	8%	7%	14%

- click-logs
  - additional training data





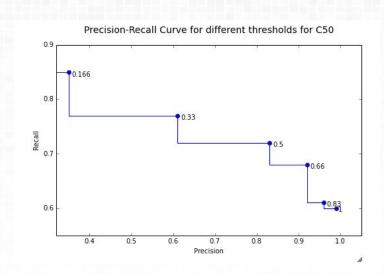
- BRAT (brat rapid annotation tool), <a href="http://brat.nlplab.org/">http://brat.nlplab.org/</a>
- highly configurable/flexible
- good for sequence labeling, relations, ...
- open-source, web-based, collaborative annotations (diffs)
- pre-labeling & correct
- extensions for IAA tooling

## **Performance Evaluation**



on CPROD (subset)

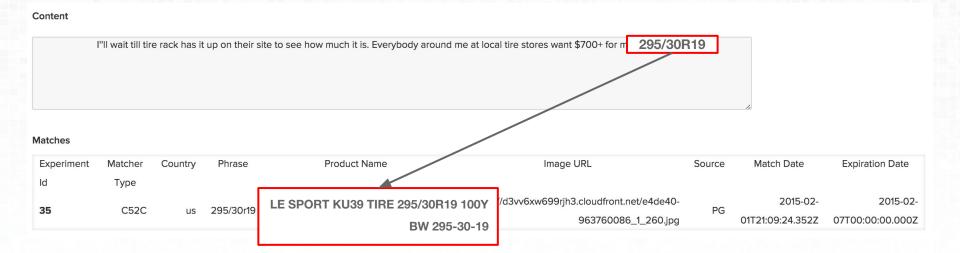
	R	Р	F
dict	75%	91%	82%
rules	69%	100%	82%
M11	76%	51%	61%
M12	78%	48%	60%
M16	72%	83%	77%
MP3	88%	98%	93%
ensemble	92%	43%	59%



- AB tests for CTR
  - sloppy models need to be in ("interesting" terms)

## **Linking Product Mentions to Offers**





## Challenges



- unstructured offers
- accessories
  - ~ 50% in Consumer Electronics
- overloaded offer titles
- nonsensical/uninformative offer titles
- diversity
  - 33 top-level industry types
  - ~50K merchants
  - ~20 feed sources
- volume
  - ~400 million offers

Polk Audio Blackstone TL1600 Speaker sys - home theater - 5.1-CH - wired

iPhone 5 Power Button Replacement

Apple iPhone 5, Space Gray, 16 GB (Unlocked), with cover, free shipping, like new

# Challenges



"iphone 5" ----

→ 650K offers!

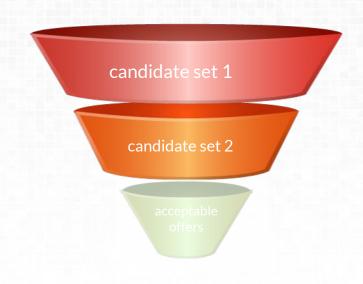
- → Armband for **iPhone 5** / 5s / 5c with a slim case (Also fits slim cases for iPhone 5, iPhone 5c, iPhone 4S and iPhone 4)
- → Genuine Porsche Martini Racing iPhone 5 iPhone 5S Case Cover
- → iPhone case with your own photo **Iphone 5** Covers
- → iPhone 5 Power Button Replacement
- → Apple iPhone 5, Space Gray, 16 GB (Unlocked), with cover, free shipping
- → .

## **Candidate Set Construction**



- semi-structured ES query
  - normalization of product mention
  - restrictions (industry type)

- business logic
  - restrictions based on user settings
  - o affiliable?



## Semantic Offer Candidate Filtering



### evaluate a tuple

cproduct mention, offer candidate>

### by these decision criteria

### 1. validity

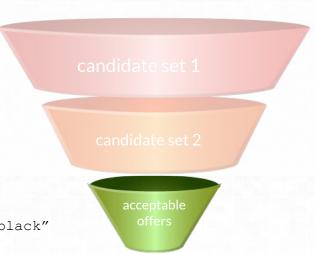
"denim jeans" // white  $\underline{\text{denim jeans}}$  jacket

### 2. centrality

"matrix amp" -> "<a href="matrix">matrix</a> mini-i 24 384 dac headphone <a href="matrix">amp</a> black"
"vudu" // "blu-ray disc player - streams <a href="vudu">vudu</a> videos"

### 3. accessory

"iphone 4" -> in-ear <u>earbuds</u> headset for <u>iphone 4</u>

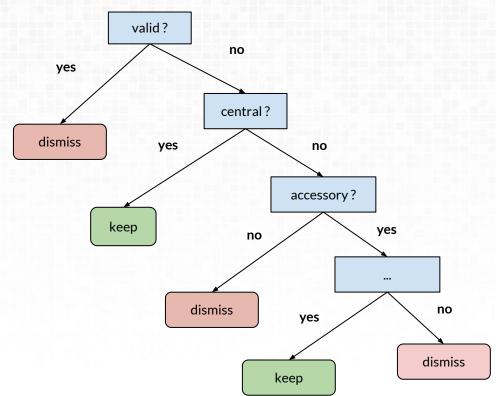


## Semantic Offer Candidate Filtering



decision logic based on these criteria

 heuristics & learning-based approaches to recognize aspects for tuples



# Manually Annotated Test Set



	[PHRASE [P]]					
1	window 7 64-bit					
	OFFER					
3	window 7 64 bit restore recovery disc					
	OFFER					
4	window 7 64 bit home premium system recovery					
	OFFER					
5	laptop window 7 64 bit system recovery disk boot dvd disc					
	OFFER					
6	usb to serial rs232 db9 adapter converter for win7 window 7 64 bit vista xp tw					
	OFFER					
7	usb to serial rs232 db9 adapter converter for win7 window 7 64 32 bit vista xp2o					
	OFFER					
8	set of 3 dell optiplex 380 sff w window 7 pro 64 bit 2gb ram 250 hd 2 93ghz					
LO	CONTEXT: WTS -> Sapphire 5770 Vapor-X VR-Zone VR-Zone () Pricewatch SGCaft Video Cards and Monitors WTS -> Sapphire 5770 Vapor-X Printer Friendly View   Emirep: 4 (100%) Infractions: 0/0 (0) WTS -> Sapphire 5770 Vapor-X fantasy87 Apr 8 rep: 10 (100%) Infractions: 0/0 (0) Thanked 1 Times in 1 Post vtngmmk Apr 9th, 1 (0) fantasy87 Apr 9th, 11, 09:24 AM # 3 thanks! ups! cwwmeng Registered User J fantasy87 Registered User Join Date: Sep 2006 Location: Singapore Posts: 159 Trade rep: 4 (100%) Infractions: 0/0 (0) fantasy87 Apr 10th in 1 Post vtngmmk Apr 10th, 11, 10:02 PM # 7 \$130 up. fantasy87 Registered User J fantasy87 Registered User Join Date: Sep 2006 Location: Singapore Posts: 159 Trat Location: Singapore ( Pasir Ris ) Posts: 133 Trade rep: 6 (100%) Infractions: 0/0 ( P45 Memory:4X2GB Kingston 800 DDR2 HDD:Western Digital Caviar Black WD5001.					
	Optical:Window 7 64-Bit Case:Cooling Master HAF 922 w/2 200mm Cooling Master R					

Google, Wikipedia, Ebay, ES	
Entity type	
□ OFFER	
○ INVALID_OFFER	
□ ○VALID_MATCH_OFFER	
O PHRASE_CENTRA	
PHRASE_NOT_CEN	NTRAL
O FINISHED	
PARTIALLY_FINISHED     PROBLEMS	
GIGNORE	
PHRASE	
PHRASE MENTION	
O PHRASE_MENTION	
OPHRASE_MENTION	
PHRASE_MENTION  Entity attributes  offer_Is_accessory	
Entity attributes	
Entity attributes	

## Manually Annotated Test Set



### • for 450 product mentions

	AU	CE	FS	avg
valid	68%	79%	86%	80%
central	21%	37%	86%	59%
accessory	79%	51%	3%	31%

### • recognition performance

- for final decision (keep/dismiss)
- ~ 70%

## Offer Feeds



- ~20 feeds (some of which come as hundreds of files)
- ~400 million offers
- pulled daily

### Offer Store



**Apache** 

v1: Tomcat & Lucene

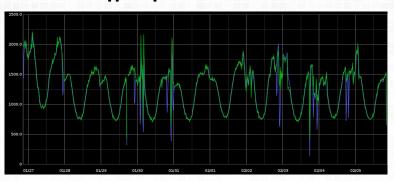
v2: Solr 3.x (x < 5) on a single box with lots of memory

v3: Solr 3.5 with (productId mod 5) sharding and ELB load balancing. Early 2013 - early 2014

### Offer Store v4: elasticsearch



### 2-10k qps peak

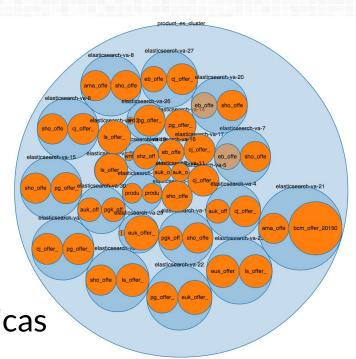




30 m3.2xlarge nodes

1-8 shards/index. no replicas

indices unioned with alias



### Raw Offer Example



```
<item data>
<item basic data>
<item unique id>B00D4INHD2</item unique id>
<item sku>B00D4INHD2</item sku>
<parent_asin>B008U04Y4Q</parent_asin>
<item upc>848447006694</item upc>
<item mpn>MH8U2AM/A</item mpn>
<item brand>Beats</item brand>
<item manufacturer>Beats Electronics, LLC</item manufacturer>
<item_name>Powerbeats by Dr. Dre In-Ear Headphone (Neon Pink)</item_name>
<item category>CE</item category>
<item_short_desc>&lt;div class="aplus"&gt; &lt;div class="leftImage" style="width: 32
ss="leftImage" style="width: 325px"> <img src="http://g-ecx.images-amazon.com/
http://g-ecx.images-amazon.com/images/G/01/aplus/detail-page/B00D4INHD2_3-sm.jpg" alt
ts.</p&gt; &lt;h4&gt;Features&lt;/h4&gt; &lt;ul&gt; &lt;li&gt;Powerful bass and hi
s for a perfectly sealed fit</li&gt; &lt;li&gt;Actual cord length: 4.3ft/1.3m&lt;/
```

#### **Normalized Offer**



```
"_index": "euk_offer_20150209_wqm6y",
"_type": "offer",
" id": "IDLkHGw8IvY",
 score": 7.1824985,
▼ "_source": {
   "offerId": "IDLkHGw8IyY",
   "productId": null,
   "displayName": "Dr Dre Powerbeats Pink",
   "name": "dr dre powerbeats pink",
   "targetUrl": "http://rover.ebay.com/rover/1/710-53481-19255-0/1?
   toolid=10029&campid=CAMPAIGNID&customid=CUSTOMID&catId=293&tvpe=2&ext=161
   "imageUrl": "http://i.ebayimg.com/00/s/NTMyWDYwMA==/z/YX4AAOSw7ThUiwOJ/$ 1.JP
   "expirationDate": "201501140000",
   "source": "EUK",
   ▼ "categories": {
       "1": "CE"
   "listingId": "161522090101",
   "feedMerchantId": null,
   "networkProductId": null,
   "merchantName": "ebay.co.uk",
   "epc": null,
   "price": 83,
   "country": "uk",
   "category": "CE",
   " timestamp": 1418848860000,
    ▼ "asin": [ ],
    ▼ "upc": [ ],
   " ttl": 2344740000
```



### Offer Import Streams



- parse (csv, xml, ...)
- extract fields
- write to elasticsearch as \_write
- write to s3 as ison
- swap to \_read when complete
- drop old index
- 3-day TTL



Akka Streams coming soon (async; backpressure)

### Offer Understanding



"Polk Audio Blackstone TL1600 Speaker sys- home theater - 5.1-CH - wired"



Brand Name (BN): Polk Audio

Product Line (PL): Blackstone

**Product (PI):** <u>TL 1600</u>

Product Category (PC): speaker sys

**Product Feature (PF):** <u>home theater</u>

Product Feature (PF): <u>5.1-CH</u>

Product Feature (PF): wired

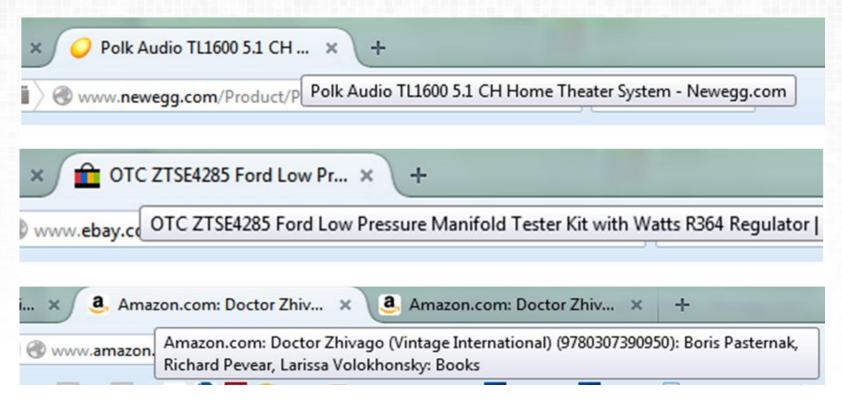




[Polk Audio]<sup>BN</sup> Blackstone<sup>PL</sup> TL1600<sup>PL</sup> [Speaker sys]<sup>PC</sup> [home theater]<sup>PF</sup> [5.1-CH]<sup>BN</sup> [wired]<sup>PF</sup>

#### >~ 1B Offers w/ Unstructured Titles





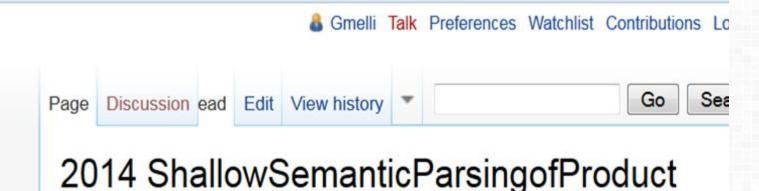
# Gabor Melli's Research Knowledge Base

Navigation

Main page Recent changes Random page

Tools

What links here



 (Melli, 2014) ⇒ Gabor Melli. (2014). "Shallow Semantic Parsing of Product Offering Titles (for Better Automatic Hyperlink Insertion) ■." In: Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining. ISBN 978-1-4503-2956-9 doi:10.1145/2623330.2623343 ☑

Subject Headings: Automated Annotation, Product Offer Title.

**Notes** 

Presentation slides: File:2014 ShallowSemanticParsingofProduct.kdd2014.pdf

O:1- - I D. .

## Manually Annotated Records



Annotated product offering title	Industry
Tripp <sup>PI</sup> LS606M <sup>PF</sup> [600 watt] <sup>PC</sup> [Line Conditioner] <sup>PF</sup> [6 Outlet] <sup>PF</sup> [120 volt] <sup>QQ</sup>	Electronics
$Samsung^{PF}$ 51- $Inch^{PF}$ 720 $p^{PF}$ 600 $Hz^{PF}$ $Plasma^{PF}$ $HDTV^{PC}$ ( $Black^{PF}$ )	Electronics
[Wheelskins] <sup>PF</sup> [Genuine Deerskin] <sup>PC</sup> [DRIVING GLOVES] <sup>PC</sup> - Tan <sup>PF</sup> ( [Size Large] <sup>PF</sup> )	Automotive
$[1990\text{-}1997]^{\texttt{PF}} \ [Mazda]^{\texttt{BN}} \ [Miata]^{\texttt{PL}} \ AxleBack^{\texttt{PL}} \ [Exhaust \ Bolt]^{\texttt{PC}} \ on^{\texttt{FT}} \ Muffler^{\texttt{PC}}$	Automotive
[Cloudy with a Chance of Meatballs] <sup>PI</sup> ( Two-Disc <sup>PF</sup> [Blu-ray / DVD Combo] <sup>PF</sup> )	Entertainment
[A Pius Man : A Holy Thriller] <sup>PI</sup> ( [The Pius Trilogy] <sup>PL</sup> ) ( [PF  Volume 1] <sup>PF</sup> )	Books
$[Amazing\ Herbs]^{\tt BN}\ [Black\ Seed]^{\tt PF}\ [Cold-Pressed\ Oil]^{\tt PC}\ -\ 32oz^{\tt PF}$	Health & Beauty
( [180 days wrty.] OF ) ZeroLemon BN LGBN Nexus PL 4 PI Juicer PL [Battery Case] PC	Mobile
$[WORLDS\ ONLY]^{\tt OF}\ NEXUS^{\tt PL}\ 4^{\tt PI}\ [BATTERY\ CASE]^{\tt PC}\ (\ LG-N4-BattCase-black^{\tt PI}\ )$	Mobile
[Kenneth Cole] <sup>BN</sup> REACTION <sup>PL</sup> men's <sup>PF</sup> Cufflinks <sup>PC</sup> , Silver <sup>PF</sup> , [One Size] <sup>PF</sup>	Fashion

token	First Char	Char Count
ViewSonic	FRSTCHR_UC	CHARCNT_9
3D	FRSTCHR_num	CHARCNT_2
HD	FRSTCHR_UC	CHARCNT_2
Video	FRSTCHR_UC	CHARCNT_5
Processor	FRSTCHR_UC	CHARCNT_9
(	FRSTCHR_oth	CHARCNT_1
VP3D1	FRSTCHR_UC	CHARCNT_5
)	FRSTCHR_oth	CHARCNT_1
Netgear	FRSTCHR_UC	CHARCNT_7
85Mbps	FRSTCHR_num	CHARCNT_6
Powerline	FRSTCHR_UC	CHARCNT_9
Network	FRSTCHR_UC	CHARCNT_7
Adapter	FRSTCHR_UC	CHARCNT_7
Kit	FRSTCHR_UC	CHARCNT_3
=-	FRSTCHR_oth	CHARCNT_1
XETB1001	FRSTCHR_UC	CHARCNT_8

Training Set					Test S	et Cat	egory				
training set	All	CE	BK	HG	CM	НО	JW	НВ	FS	AU	AE
All	57.7	64.3	41.2	54.6	58.3	48.0	60.3	51.6	55.0	46.7	48.8
CE	47.1	63.3	6.6	46.6	48.7	34.8	35.8	34.6	36.6	42.7	27.2
BK	11.1	4.0	82.0	1 2.8	4.2	2.6	2.1	3.5	3.8	2.8	37.5
HG	34.9	41.4	7.5	50.0	28.9	32.8	36.1	32.3	32.3	32.9	28.1
CM	41.2	48.9	9.1	39.7	56.5	29.8	32.5	33.1	31.0	35.7	31.5
НО	27.3	30.4	4.1	37.3	22.3	43.5	28.7	29.6	25.7	27.7	18.9
JW	27.1	30.2	7.2	31.8	19.5	23.9	56.2	29.3	31.9	26.1	15.9
HB	27.9	30.9	4.7	35.1	23.3	29.2	27.3	45.8	28.1	24.3	23.6
FS	27.0	32.0	5.8	31.7	22.1	23.1	32.1	27.6	52.3	25.0	27.8
AU	31.1	40.0	4.8	33.9	24.8	26.8	31.4	25.0	26.7	38.2	19.2
AE	13.3	9.0	59.2	6.3	7.5	7.4	4.3	9.7	11.7	4.6	61.3

.

\_

Training Set					Test S	Set Cat	egory				
training set	All	CE	вк	HG	CM	но	JW	НВ	FS	AU	AE
All	57.7	64.3	41.2	54.6	58.3	48.0	60.3	51.6	55.0	46.7	48.8
CE	47.1	63.3	6.6	46.6	48.7	34.8	35.8	34.6	36.6	42.7	27.2
BK	11.1	4.0	82.0	2.8	4.2	2.6	2.1	3.5	3.8	2.8	37.5
HG	34.9	41.4	7.5	50.0	28.9	32.8	36.1	32.3	32.3	32.9	28.1
CM	41.2	48.9	9.1	39.7	56.5	29.8	32.5	33.1	31.0	35.7	31.5
НО	27.3	30.4	4.1	37.3	22.3	43.5	28.7	29.6	25.7	27.7	18.9
JW	27.1	30.2	7.2	31.8	19.5	23.9	56.2	29.3	31.9	26.1	15.9
HB	27.9	30.9	4.7	35.1	23.3	29.2	27.3	45.8	28.1	24.3	23.6
FS	27.0	32.0	5.8	31.7	22.1	23.1	32.1	27.6	52.3	25.0	27.8
AU	31.1	40.0	4.8	33.9	24.8	26.8	31.4	25.0	26.7	38.2	19.2
AE	13.3	9.0	59.2	6.3	7.5	7.4	4.3	9.7	11.7	4.6	61.3
BK&AE	12.9	6.2	65.2	3.6	6.0	4.0	3.0	5.2	4.9	3.3	65.3
not BK AE	55.8	64.6	18.1	54.3	59.0	49.3	58.4	52.1	56.3	49.2	34.9

	graphics card	directional antenna		am radio	messeng		
9/1/2013	Caru 6	antenna	bag		er bag	mouse	printer
		4			8		
9/2/2013	10						
9/3/2013	3		2				
9/4/2013	7	S.			0		
9/5/2013	5	3.			2	51	
9/6/2013	4						
9/7/2013	6						
9/8/2013	7						
9/9/2013	9			1			
9/10/2013	9				9	1	
9/11/2013	6						
9/12/2013							
9/13/2013	4	2					
9/14/2013	8						
9/15/2013	7	2	2				1
9/16/2013	5	6		1			
9/17/2013	5		3				1
9/18/2013	16	6	6	3	3		2
9/19/2013	28	2	3	Ĭ	3		2
	28	4	5		2	1	
9/20/2013		A CONTRACTOR OF THE CONTRACTOR		1		1.00	
9/21/2013	33	10	3	1	1	2	1
9/22/2013	35	()		1	2	4	
9/23/2013	16	6	6	3	1	3	

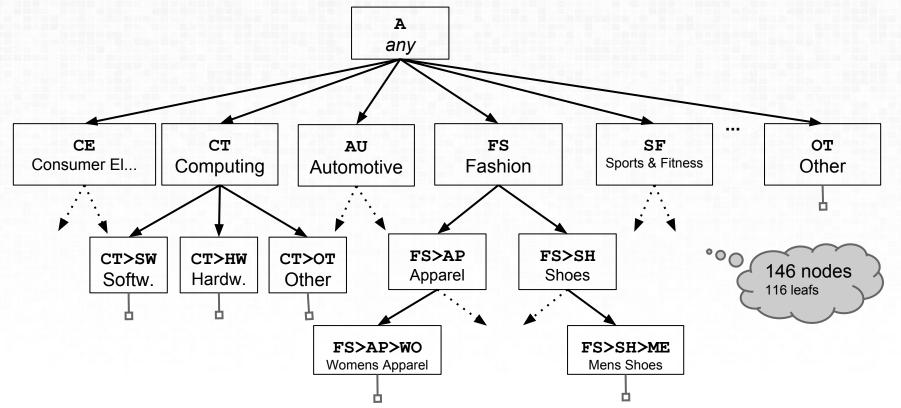
# Active Iterative (Re) Training



- Add examples with low confidence predictions.
- Preannotate new examples so that the annotator can correct:
  - [token1 token2]<sup>PF</sup> instead of [token1]<sup>PF</sup> [token2]<sup>PF</sup>
  - PF instead of PC

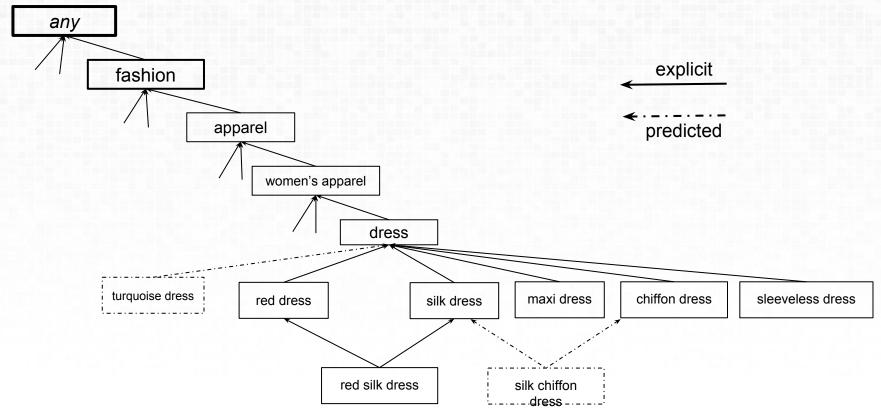
## **Product Category Hierarchy**





### Product Category Future | Subsumption (isA) Lattice





#### Concept Discovery (word2vec)



Cpb-400

Cpb-400b

Cpb-400c

Cpb-403b

Cph-485

Cph-489

Cph-496

Cph-514

Cph-516b

Cp1-536

Cpp-529z

Ct-11

Cx7525

Cxai5198

Cxai5698

..



Pwrc51

Pwrc61

Pwrc62

R-5502-W

R-5650-s

R-5650-w

R5800

R-5800-w

R5800wii

R800

Rb-41

Rb-51

Rb-61

Rc-52

Rc60i

...





#### Thank You!



# Interested in working for VigLink?

#### See open positions at

http://www.viglink.com/about/careers/openings/