# Connecting The Probable Dots

Disambiguating strategies for entity oriented search and natural language understanding

## Dawn Anderson

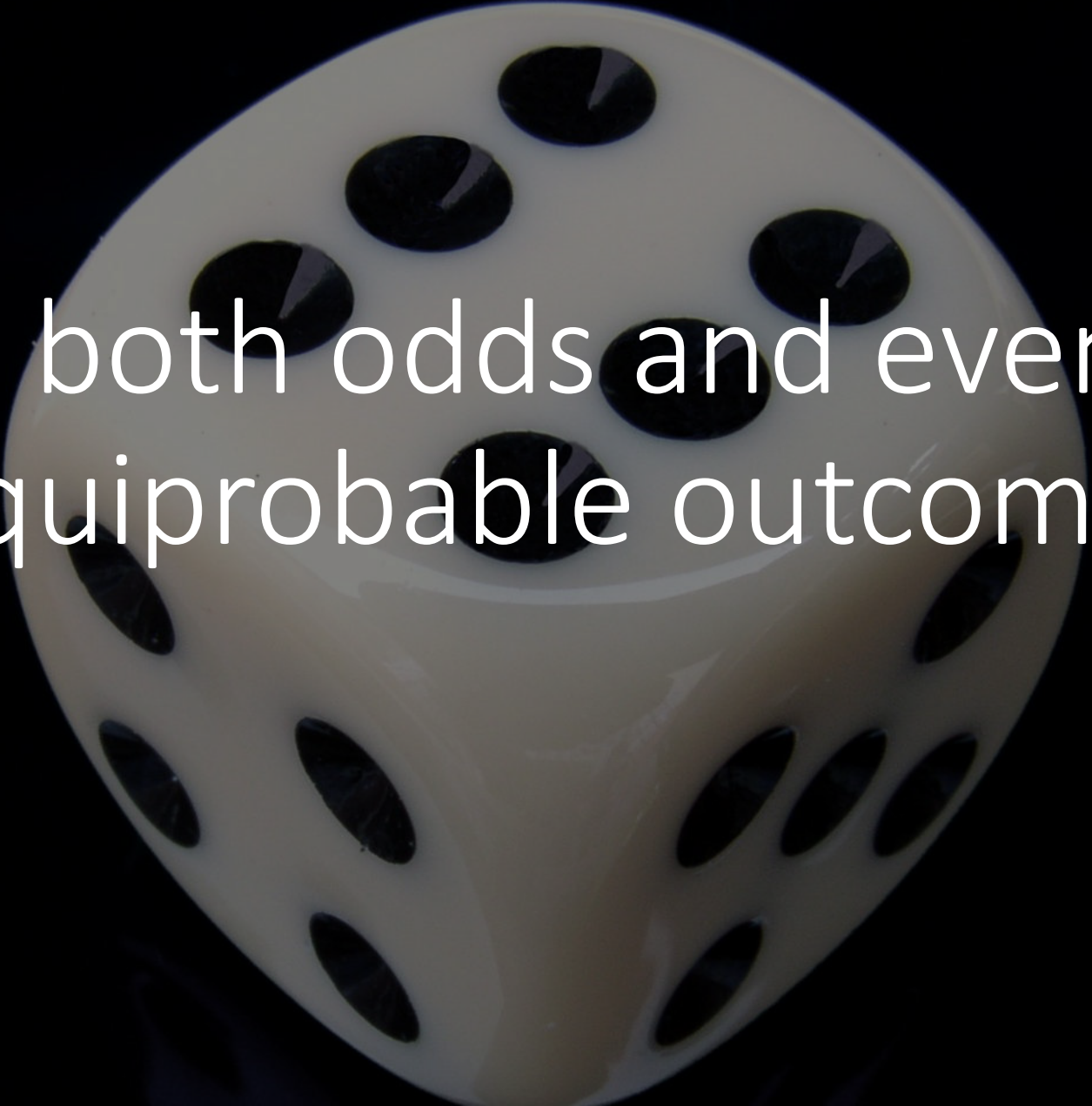@dawnieando                    #FOS20

The probability of rolling an even number or an odd number is 1:2

This is an example of 'equiprobability'

Since both odds and evens are equiprobable outcomes

"Assigns equal probabilities to outcomes when they are judged to be equipossible" (Wikipedia)

In URLs, content, locations & entities equiprobability can be problematic

When one of a number of pluralities equally meets an information need

One must beat the other

THERE CAN BE ONLY ONE

This could be duplicate content or 'ambiguity in entities'

Equiprobability is like 'See saw SEO'

# LaPlace's Principle of Indifference

la places principle of indifference

All    Images    Videos    News    Shopping    More      Settings    Tools

About 8,260,000 results (0.53 seconds)

## Principle of indifference philosophy

The **principle of indifference** states that in the absence of any relevant evidence, agents should distribute their credence (or 'degrees of belief') equally among all the possible outcomes under consideration. In Bayesian probability, this is the simplest non-informative prior.

en.wikipedia.org › wiki › Principle_of_indifference

**Principle of indifference - Wikipedia**

Search for: Principle of indifference philosophy

About Featured Snippets      Feedback

# Where there is equi-possibility there is a need for further confirmations

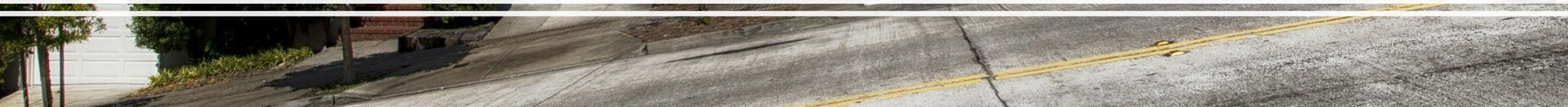Further 'confirmations' beyond the same content or surface form must be sought

Where two or more domain assets, surface forms, or entity determinations are considered equal a representative is 'the chosen one' (canonical)
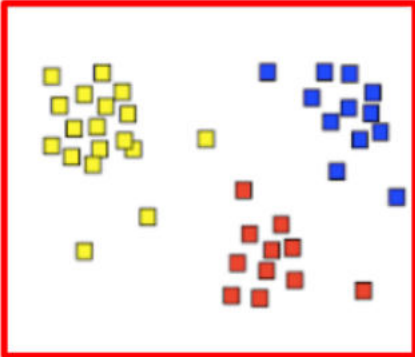
Nearest Neighbours

All the pluralities are clustered together and one is picked

clusters

Q All  Images  News  Shopping  Videos  More  Settings  Tools

About 492,000,000 results (0.99 seconds)

**Cluster** analysis or **clustering** is the task of grouping a set of objects in such a way that objects in the same group (called a **cluster**) are more similar (in some sense) to each other than to those in other groups (**clusters**).

Cluster analysis - Wikipedia
https://en.wikipedia.org › wiki › Cluster_analysis

Sometimes 'bits' are picked from multiple 'equiprobable outcomes'

**Near Duplicate Content Can Cause Google to Choose Wrong Snippet**

MARCH 14, 2017 AT 6:30 AM PST BY JENNIFER SLEGG

**TheSEMPost**

Near Duplicate Content Can Cause Google to Choose Wrong Snippet

Today's Topic: Disambiguating equiprobability for improved search performance

WHO IS DAWN ANDERSON

linkedin.com/in/msdawnanderson

@dawnieando

dawn.anderson@ @move-it-marketing.co.uk

11+ years SEO & Digital Marketing Consultant & Pracademic:

Manchester Metropolitan University

Move IT
Digital Marketing Solutions

Contributor:

Search Engine Land

Smart Insights

TheSEMPost

SEJ Search Engine Journal

Speaker & Trainer:

SMX LONDON UK

SMX WEST SAN JOSE

MozCon

brightonSEO.

unGAggED

PUBCON

SAScon
SEARCH ANALYTICS SOCIAL MEDIA

STATE OF SEARCH CONFERENCE

GENERAL ASSEMBLY

Squared Online

# WHO IS DAWN ANDERSON

linkedin.com/in/msdawnanderson

@dawnieando

dawn.anderson@ @move-it-marketing.co.uk

move-it-marketing.co.uk

Meet Bert & Tedward

In an ideal world 'precision' and 'recall' would be perfect... but... it's NOT

So search engines must also work off probability determination

Machine learning also now powers some probability determination

# Google Uses Machine Learning for Crawling, Indexing & Ranking

# Using 'The Law of Large Numbers'

Progressive learning & increasingly 'educated' guesses

Sometimes 'probability' predictions can be unpredictable

On the 'long tail' there's often very little, or nothing in it

Low to no page-rank pages (but that is also now just one of very many things

There will be lots of contributing factors now

## 1 - PAGE IMPORTANCE CONTRIBUTORS??

- Location in Site (e.g. home page more important than parameter 3 level output)
- PageRank
- Page type / file type ('about us' e.g. less important)
- Inclusion in XML sitemap (if others are excluded)
- Internal PageRank
- Internal Backlinks
- In-site Anchor Text Consistency
- Relevance (content, anchors and elements) to a topic (Similarity Importance)
- Directives from in-page robot and robots.txt management
- Parent quality brushes off on child page quality

**IMPORTANT PARENTS LIKELY SEEN TO HAVE IMPORTANT CHILD PAGES**



FIG. 4

See-saw SEO

Then the canonical balance Tips From Equiprobability Slightly

Your ranking flux might well be transferring equiprobability

Several types of equiprobability ambiguitity may cause your project to perform less well overall in search

# Different types of ambiguity can impact SEO

Exact duplicate content

Near duplicate content

Natural language ambiguity

Generational cruft based ambiguity

Machine learned ambiguity 'lag'

Dot-to-dot ambiguity

Semantic heterogeneity

Location based ambiguity

# Exact duplicate content

Google probably chooses the one with highest 'probability'

Most linked to internally

Most linked to externally

Included in XML sitemap

HTTPS

Prettier URLs

# Maybe one or several of many 'importance' factors

---

## 1 - PAGE IMPORTANCE CONTRIBUTORS??

- Location in Site (e.g. home page more important than parameter 3 level output)
- PageRank
- Page type / file type ('about us' e.g. less important)
- Inclusion in XML sitemap (if others are excluded)
- Internal PageRank
- Internal Backlinks
- In-site Anchor Text Consistency
- Relevance (content, anchors and elements) to a topic (Similarity Importance)
- Directives from in-page robot and robots.txt management
- Parent quality brushes off on child page quality

**IMPORTANT PARENTS LIKELY SEEN TO HAVE IMPORTANT CHILD PAGES**



FIG. 4

Google works this out VERY quickly

# Near-duplicate content

Easy(ish)

It starts off well enough... but over time

Google realizes these pages are 'mostly' the same as others

# Some symptoms

Many matching shingles

Quilting

'Borrowed' content at scale

Feeds

Big Boilerplate / little main value

Data driven sites with little value

Sometimes these are sites trying to make a URL footprint well beyond their contribution to a positive 'Network Effect'

After some crawling & sampling Google puts the pieces of the puzzle together

And begins to exclude pages from the index

Exclusion still happens... A LOT

# Particularly on behemoth sites

# I ran a little Twitter competition with only empathy as the prize

**Dawn Anderson**
@dawnieando

Competition time:

Who's got the highest number of 'Excluded' in any GSC Coverage Report.

Show me the screenshots.

No prize for winning.  Just the glory and empathy for the work ahead ;P ;P

11:00 AM · Jan 27, 2020 · Twitter Web App

I would hazard a guess…The majority of the largest websites in the world have a LOT of near-duplicate 'transactional' pages

Probably 20% of the URLs satisfy 80% of the demand & contribute well to 'The Network Effect'

Since search engines probably work with positive 'impact' value considered

Google will likely use
mostly 'Sampling' on those
heavy 'excluded URL' sites

Looking for small samples of content & URL patterns… and… 'Discovered, not crawled' will be high

'Discovered, not crawled' in GSC is Google saying...

They know what's down that site section or URL parameter path

The verdict is in _____


THEY KNOW IT'S A LOAD OF POOP

# Those URLs will be 'starved' of crawl

However many times you submit and inspect

You have your work cut out once Google puts the pieces of that puzzle together

# Natural language ambiguity

**For machines words** are problematic.  Ambiguous… polysemous… synonymous

In spoken word it is even worse because of homophones and prosody

Like "four candles" and "fork handles"

Many words have multiple meanings.

Like "like" can be 5 possible parts of speech (POS)

# Word's context helps enormously

And disambiguating words is getting easy(ier)

"You shall know a word by the company it keeps" (Firth, 1957)

Words that are related live near each other in mathematical spaces

There are HUGE leaps forward in natural language understanding now using machine learning

Accelerated mostly by Google's BERT (Bi-Directional Encoder Representations from Transformers)

**Co-occurrence in content helps A LOT**

In the page itself

In the interconnected pages

In the subcategorisation

In the site sections

In the external relationships

In the domain ontology as a whole

# Relatedness

First level relatedness

Second level relatedness

Informational content adds to the contextual 2nd level relatedness of transactional content

Informational content

Transactional content

If you can help with 'informational needs' you can 'probably' help with transactional needs

Add value in rich informational needs

Contextual value passes throughout the whole site

# Adding Contextual Value Does Not Mean

| | |
|---|---|
| **Adding** | Adding loads of content below ecommerce pages |
| **Adding** | Adding LSI keywords in ecommerce pages |

High topical relatedness & context vectors

Internal linking

Not spam

Utilise conceptual 'nearest neighbours' well

Identify content which
is 'very' closely related

Merge content which is 'too' conceptually very similar with no separate demand

# Utilise 'Overflow SEO' as demand & content naturally grows

Google may realise you can 'probably' help with transactional queries

# Dot-to-dot ambiguity

Who knows what this dot to dot puzzle is creating?

Absolutely no-one

If you chop away or migrate half of your website then equiprobability will likely be a problem

A half drawn website is not the same as a fully drawn previous website?

Informational content adds to the contextual 2nd level relatedness of transactional content

Informational content

Transactional content

Although you may rank 'marginally' better for fewer things

Your domain 'topic' became overall more about some topics & less about other topics

# Prune, Merge, Improve, Archive

## 01
Prune – Only ballast

## 02
Merge – Highly semantically related content

## 03
Improve – Evergreen pages worth the effort

## 04
Archive – Older pieces of temporal content (like a library)

Avoid removing or hiding user generated content (except spam)

Google puts the pieces of a half-site puzzle together (often with bad results)

# Cruft based ambiguity

Generational cruft contributes to equiprobability issues

Just a few examples... legacy is a problem

301 redirect chains

Inconsistent 301s

No 301s at all

Canonical points to a 301

Canonical points to a no-index

No-index points to a canonical

# 301 Should Mean The Resource Moved

You're supposed to be telling search engines where the words (tokens), topics, concepts & entities went

You're saying things moved so they can be re-filed



Document 1

The bright blue butterfly hangs on the breeze.

Document 2

It's best to forget the great sky and to retire from every wind.

Document 3

Under blue sky, in bright sunlight, one need not search around.

**Stopword list**

a
and
around
every
for
from
in
is
it
not
on
one
the
to
under

**Inverted index**

| ID | Term | Document |
|----|----------|----------|
| 1 | best | 2 |
| 2 | blue | 1, 3 |
| 3 | bright | 1, 3 |
| 4 | butterfly | 1 |
| 5 | breeze | 1 |
| 6 | forget | 2 |
| 7 | great | 2 |
| 8 | hangs | 1 |
| 9 | need | 3 |
| 10 | retire | 2 |
| 11 | search | 3 |
| 12 | sky | 2, 3 |
| 13 | wind | 2 |

But often there is little or no match at all

Borderline or true 'soft 404'


MOVE ALONG
NOTHING TO SEE HERE
imgflip.com

Also, any legacy canonicals are now defunct by the redirect status

Cruft & migrations are like a game of rugby as the ball (signals) is passed between URLs

# A sign of this is a temporarily wrong target ranking

One URL is going up.  One is going down as signals pass from one to the other over time – See-saw SEO

Google puts the pieces of the puzzle together... eventually (or never)

Machine learned ambiguity 'lag'

# ALL of the above ==

# Crawl Budget Woes / Learned Quality Patterns

# Your quality & crawl will be 'machine learned' by 'Large Numbers' over time

Probably you'll struggle to get enough crawl for Google to catch up

Search engines realise there is no 'demand' for your poopy pages'

You've probably got URLs which have not been crawled for years

Your performance will be based on what is indexed

# Fear not... Small Wordpress site

You need to lure a Grumpy Googlebot with tasty quality content morsels

You MUST change the ratio of poor quality to high quality & add more value OVERALL

And contribute positively to 'The Network Effect'

You need to substantively improve the quality of your pages, regain the crawl and get more demand

Watch for patterns (clues) in GSC Coverage & take a demand-driven approach

# In transactional pages identify valuable content, features & attributes your audience wants

Pro-tip: Adding moooaaarr words to transactional pages does not always add much value

Plus… There is often a mathematically natural ordered pattern to things ranked by frequency

Like word frequencies many types of 'ordered' popularities will have a Zipfian Distribution

Where the frequency of x is inversely proportional to its frequency table rank – 1/n

Zipfian Distribution occurs in many other rankings unrelated to language

Population of cities in a country

Corporation sizes

Income rankings

No. of people watching same TV channel

# Popularity & Zipfian Distributions

Is it really popular in 'real life'?

Does it follow 'The Network Effect'?

Or are you manipulating it for things you want to rank for (which aren't really popular at all)?

# Bow tie of the web & strongly connected components

Building genuine topical hubs (Hub & spoke)

Google puts the pieces of the puzzle together slowly but increases with quality improvements

Equipossibility goes well beyond duplicate content too

# Semantic heterogeneity

# The Web of Document Vs The Web of Data

# Applies mostly to entity-oriented search

# Mis-matching data types or equivalencies in data tables from same or other domains

# Linked Data has been around for a long time



The **inventor** of the World Wide Web and the creator and advocate of the Semantic Web and **Linked Data**, Sir Tim Berners-Lee, laid down the four design principles of **Linked Data** as early as in 2006.

www.ontotext.com › knowledgehub › fundamentals › linked-data-linke...

What are Linked Data and Linked Open Data? - Ontotext

Link to more data

URI

Resource on the web

HTTP access

# There have historically been several ways to implement linked data

Several linked-data types & web sources come together causing equipossibility

# Some Types of Semantic Web Technologies & Their Markups

RDF

SPARQL

OWL

SKOS

RDFa

JSON-LD

Microdata

# Implementation Inconsistencies Prevail (ed)

The Knowledge Graph & Its Data Sources

The majority of Google's entity types reference the Wikipedia URI or the Knowledge Graph MID

# But Now Anything Structured Might Be Used



**GOOGLE MAY USE ENTITY EXTRACTIONS, ENTITY CLASSES, ENTITY PROPERTIES, AND ASSOCIATION SCORES FROM PAGES TO BUILD KNOWLEDGE GRAPHS**

When Google introduce the Knowledge Graph in 2012, it told us that it was going to start focusing upon things and not strings. That process is maturing, and we have a chance to watch Google learn how to start crawling the Web to mine data and engage in entity extractions, instead of mining web information such as pages and links. As I wrote recently on Twitter about this:

> **Bill Slawski** ⚓
> @bill_slawski
>
> In web crawling, a node is a page, and an edge is a link between pages; in data crawling, a node is an entity, and an edge is a relationship between entities. It's an evolution in thinking about the web.
>
> ♡ 274   2:58 PM - Feb 10, 2019 · Carlsbad, CA   ⓘ
>
> ○ 105 people are talking about this                    >

A recently granted Google patent tells us about how the search engine may perform entity extractions from web pages, and store information about them. This goes beyond using knowledge bases as sources of information about entities, and moves on to finding more than what may be available in such sources, by looking at textual passages on web pages. The problem that this patent is intended to solve is described in this early line from the patent:

*Conventional knowledge bases, however, can fail to provide up-to-date or reliable information regarding*

# We already know conversation search fills gaps via web information

The web is a 'Data Lake'

# Public Datasets (e.g. Kaggle) > 27k datasets

# Enter Data Search by Google AI

# Disambiguating Data Sets (is hard)

# Probability Predictions Can Be VERY Wrong

# Example: Philosophers Dates of Birth different in Wikidata & DBPedia

**Ivo Velitchkov** @kvistgaard · Jan 27

@dawnieando Here is a query w.wiki/GDq giving 50 philosophers for which @dbpedia and @wikidata state different date of birth.

You may change the DBpedia type for other occupations and the limit if you want different number of results.

#SPARQL

| | | |
|---|---|---|
| Hector Zagal | 1966-06-06 | 1952-06-06 |
| John Gardner | 1965-03-25 | 1965-03-23 |
| Paul Copan | 1962-09-26 | 1962-09-20 |
| Muhammad Tahir | 1962-03-21 | 1962-01-06 |
| Matthew Kramer | 1959-06-09 | 1959-01-01 |
| Mohsen Kadivar | 1959-06-08 | 1959-06-07 |
| Norbert Schmitt | 1956-01-23 | 1956-01-01 |
| Peter M. Haas | 1955-01-23 | 1955-01-25 |
| Michael Sandel | 1953-03-05 | 1953-05-03 |
| Javed Ahmad Ghamidi | 1952-04-07 | 1951-04-18 |
| Hitoshi Nagai | 1951-11-10 | 1951-01-01 |

# I asked for some examples on Twitter

**Dawn Anderson**
@dawnieando

Hey folks. Can anyone provide me with examples of inconsistencies in knowledge graphs (e.g. some data taken from one place and some taken from another place). e.g. photos from one place and content from another place). I know these are out there but please share examples

11:36 AM · Jan 26, 2020 · Twitter Web App

# The community responded

# 'Mashed-up' Knowledge Graphs (Two Don Cherry's and his dog's name - 'Blue'

---



**Aaron Bradley**
@aaranged

Disambiguation is hard. The knowledge panel is a mashup of Don Cherry, the jazz trumpeter, and Don Cherry the (former) CBC hockey commentator. Keen eyes may also have correctly assessed that that's not a picture of the dog, Blue (it's Don Cherry, "big band singer and golfer").

don cherry's dog

Q All  Images  Maps  News  Shopping  More    Settings  Tools

About 943,000 results (0.70 seconds)

Don Cherry / Dogs

## Blue

**Blue** (Don Cherry's dog) **Blue** was a white English Bull Terrier owned by hockey commentator Don Cherry. **Blue** was reportedly a gift from the members of the Boston Bruins when Cherry was their head coach from 1974 to 1979. The original **Blue**, who died in 1989, was a female.

Blue (Don Cherry's dog) - Wikipedia
https://en.wikipedia.org › wiki › Blue_(Don_Cherry's_dog)

### Don Cherry
Jazz trumpeter

**Born:** November 18, 1936, Kingston
**Died:** October 19, 1995, Málaga, Spain
**Spouse:** Luba Cherry (m. 1999), Rosemarie Cherry (m. 1957–1997)
**Teams coached:** Colorado Rockies (Head coach, 1979–1980), Boston Bruins (Head coach, 1974–1979)
**Children:** Eagle-Eye Cherry, Christian Jon Cherry, David Ornette Cherry, Tim Cherry, Cindy Cherry

Movies and TV shows    View 10+ more

Image Search Got Don Cherry's Dog Right

# Wrong price included in images from other sites

US Price Info Showing in UK 'People Also Ask'

# Edgar's Were Comical

# Van Diesel – Death – To Be Advised??



**Vin Diesel**
Actor

© Getty Images

**Vin Diesel**, born as Mark Vincent, is an American actor, producer, director, and screenwriter. He came to prominence in the late 1990s, and first became known for appearing in Steven Spielbergs Saving Private Ryan... wikipedia.org

**Born:** July 18, 1967, New York City, New York, USA

**Died:** January 30, 2014, TBA

**Height:** 5' 11" (1.82m)

**Partner:** Paloma Jiménez (2008-2014)

**Parents:** Irving Vincent, Delora Vincent

Images from Direct Competitors

Mihai Sterian
@MihaiSterian

Replying to @dawnieando

Link and content from Moqups, images from competitors.

3:29 PM · Jan 29, 2020 · Twitter Web App

# Wikipedia Made Their Point Well

"Simply saying that a horse has five legs doesn't make it true – calling a horse's tail a leg does not make it one." (Wikipedia, 2019)

A six legged horse?

Somes Reputable Knowledge Repositories are Just Plain 'Wrong'

Google Works To Fix How Many Legs Horses & Snakes Have

Apr 30, 2019 • 8:10 am | 💬 (3)
by Barry Schwartz 🐦 | Filed Under Google Search Engine

Equivalency between data sources is sought

# Efforts Are Underway to Resolve

Since 2011 major search engines have been focusing on Schema commonly

# Google Deprecating Data-vocabulary Schema

Google Sending data-vocabulary.org Schema Deprecation Notices

Jan 22, 2020 • 7:11 am | 💬 (1)
by Barry Schwartz 🕊 | Filed Under Google Search Engine Optimization



Google announced yesterday that starting on April 6, 2020 it will no longer support data-vocabulary.org schema. Hours later, Google began sending notifications via Google Search Console about this happening.

# Your Breadcrumbs Markup May Be Impacted

**Google** Search Console

## Breadcrumbs issues detected on https://www.

To the owner of https://www.          I/:

Search Console has identified that your site is affected by 1 Breadcrumbs issues:

**Top Warnings**

Warnings are suggestions for improvement. Some warnings can affect your appearance on Search; some might be reclassified as errors in the future. The following warnings were found on your site:

data-vocabulary.org schema deprecated

We recommend that you fix these issues when possible to enable the best experience and coverage in Google Search.

**Fix Breadcrumbs issues**

# Entities often have many surface forms

# Schema:SameAs and OWL:SameAs

# SameAs Seems To Be Effective

BE consistent in ALL mentions of your brand name & ALL THINGS GENERALLY

Remember: Consistency is one of the Kings

Utilise 'known' (popular) public datasets & knowledge repositories

Avoid conflict with 'known' (popular) public datasets & 'knowledge repositories'

Realise that well known knowledge repositories may be wrong

But they may be conceptually still be 'right'

They'll likely win in an equipossibility face-off

Leave no room for ambiguity where possible

Google puts the pieces of the puzzle together often badly

Location based ambiguity

# Semantic Heterogeneity is a problem locally

Entity Address (location) is a VERY Big Problem with entity-oriented search

street number

Locality - city or town

street/route name, if detected

postal code, if detected

country, if detected

broad_region - administrative area, such as the state, if detected

Narrow_region - smaller administrative area, such as county, if detected

Bristol is BOTH a county and a city

# In Scotland and Wales it's worse

Cardiff is in South Glamorgan NOT Glamorgan

Google Does Not Understand the Word 'Historic'

The passage of time adds new meaning to queries sometimes too

Older humans probably keep ambiguous county + town combinations alive

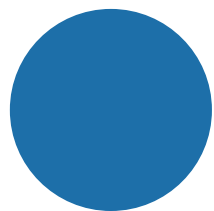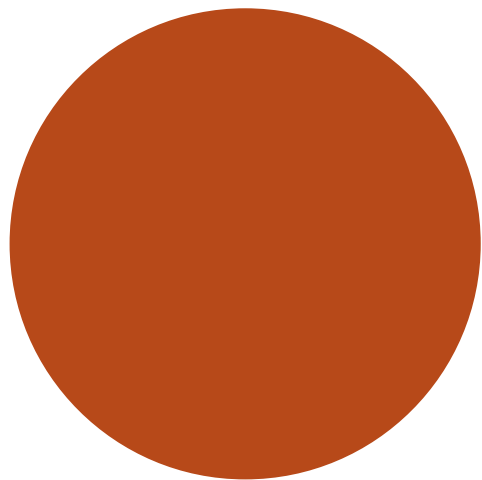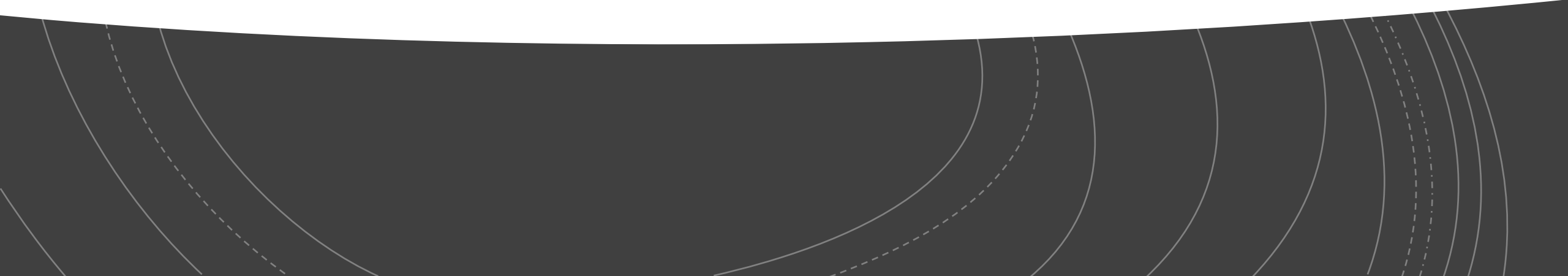# Maximise the 'Local' Opportunity

In location focused pages adding a postcode could be much more valuable than adding lots of words

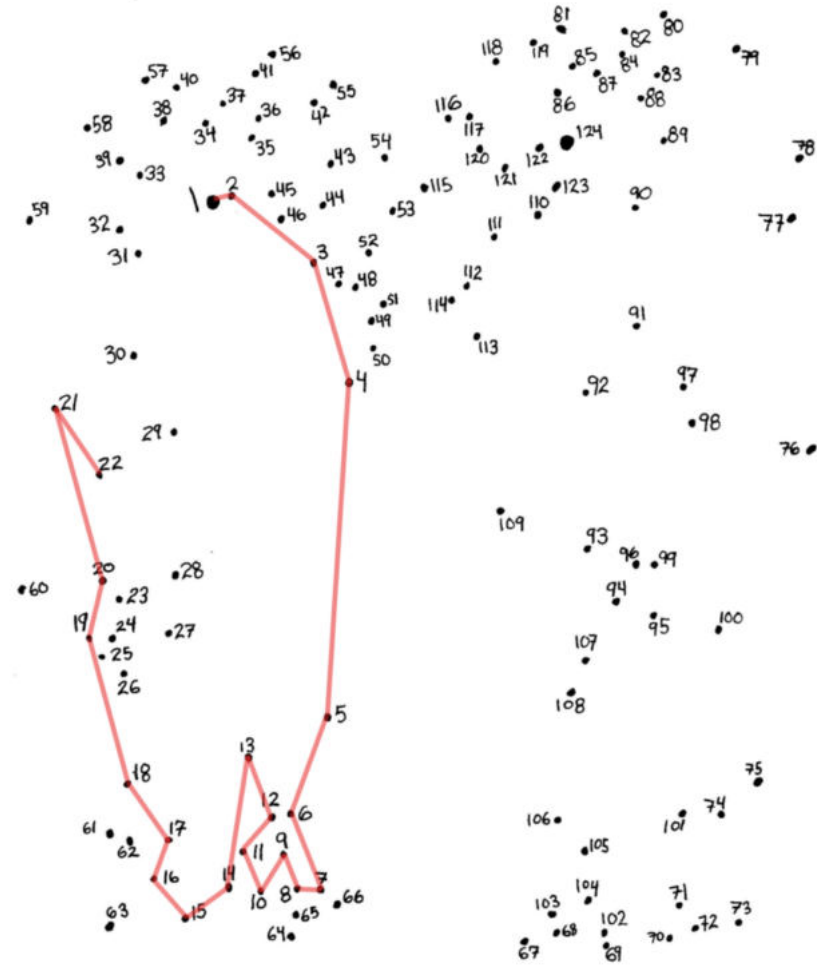# Contextual relatedness in local spaces is proximity-based

# AreaServed Schema & Internal linking by Proximity (e.g. lat & long) is powerful
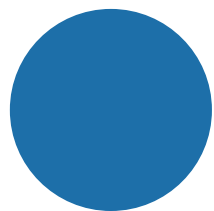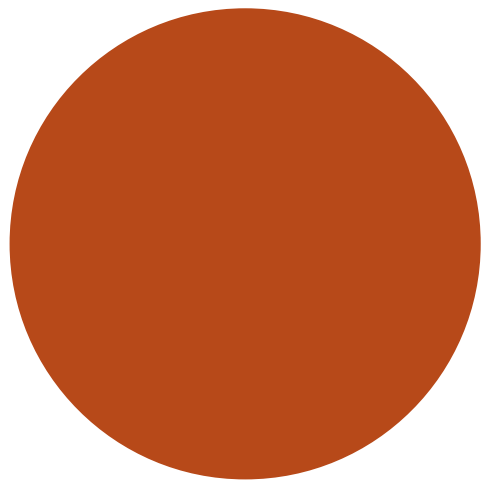
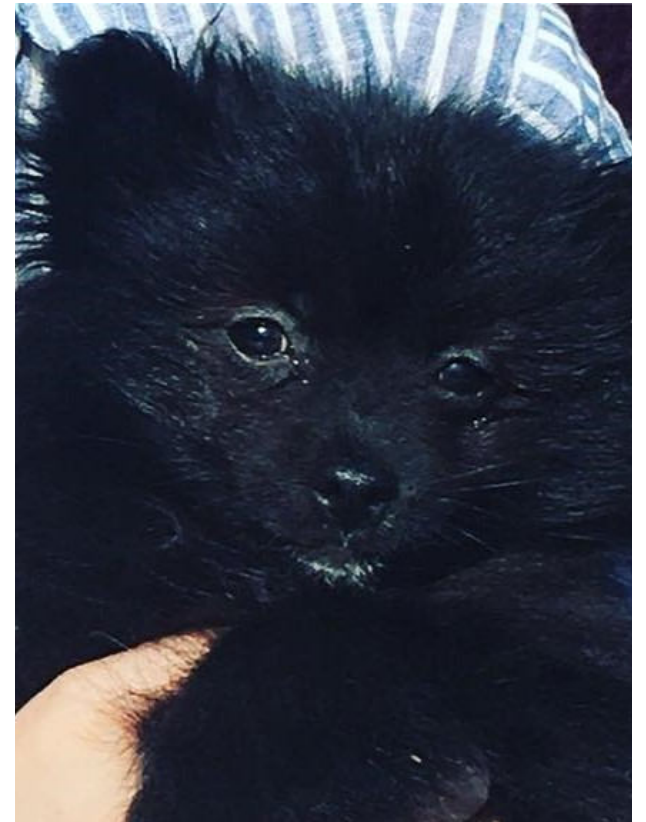Google sometimes puts the pieces of the puzzle together sporadically

Align ALL the dots

# Most Importantly - Be Consistent

# Keep in Touch

- @dawnieando
- @BeBertey

# Some Take Aways

| | |
|---|---|
| **Go** | • Watch for 'impact' patterns to maximise upon in GSC 'coverage' reports |
| **Optimise** | • Be careful you don't prune away second level relatedness signals |
| **Map** | • Consider Pareto's Law when identifying issues to fix in GSC coverage |
| **Optimise** | • Check for 'wrong page ranking' due to cruft & lag passing signals |
| **Get** | • On location sites consider internal linking on proximity |
| **Identify** | • Add value to improve on 'The Network Effect' in attributes & informational content |
| **Use** | • Make it obvious you can 'probably' help with transactional needs through other content |