

# Exploring Word2vec in Scala

Gary Sieling  
@garysieling  
Wingspan, an IQVIA Company

Jan 11, 2018  
PHASE

# FindLectures.com: A case study on natural language search

Features

Suitable for Work: 199,118

Video: 175,199

Audio: 17,084

Closed Captions: 16,234

Category

Technology: 56,723

Education: 54,324

News and Politics: 18,904

Fine Arts: 13,167

Entertainment: 11,305

+ Show More

Type

Conference: 45,606

Academic Lecture: 21,741

Historical Speech: 4,964

Interview: 4,425

Documentary: 3,465

Workshop: 805

Collection

Lanyrd: 25,962

Hacker News: 24,088

Confreaks.tv: 9,534

DPLA: 7,236

Oxford University: 6,083

+ Show More

Speaker

John Avery Lomax: 616

Ruby Terrill Lomax: 592

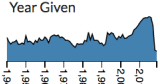
Bill Clinton: 585

Ronald Reagan: 507

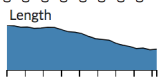
Mary Hufford: 462

+ Show More

Year Given



Length



Prison Reform: Alternatives to Mass Incarceration with Steven Raphael (29 minutes) ⓘ

Ian Bostridge and Jeremy Denk on Winterreise: Anatomy of an Obsession (60 minutes) ⓘ

Re-seeing the Unseen - Advances in Ophthalmology (15 minutes) ⓘ

Slavery, Ships and Sickness ▶

Speaker: Professor Stuart Anderson

Given On: Monday, 24 October 2011, 1:00PM

The Goat Rodeo Sessions (Yo-Yo Ma, Stuart Duncan, Edgar Meyer, Chris Thile) | Musicians At Google (34 minutes) ⓘ

Is the world flat? ▶

Speaker: Christopher Dye

Given On: Thursday, 25 October 2007, 12:00AM

Karen E. Willcox (58 minutes) ⓘ

Eva Kor: "Surviving the Angel of Death: The True Story of a Mengele Twin in Auschwitz" (69 minutes) ⓘ

Speaker: Eva Kor

Thomas L Magnanti (2.2 hours) ⓘ

Qualcomm Thinkabit Lab Presents: Programming Servos (12 minutes) ⓘ

Brad Bell and Jane Espenson - Inventing Television: How Husbands Fully Realizes the Promise (41 minutes) ⓘ

Chris Kimball: "The Science of Good Cooking" | Talks at Google (54 minutes) ⓘ

Speaker: Chris Kimball

Handling Emotion (4 minutes) ⓘ

Ethiopia 1969 Reel 5 of 65 (17 minutes) ⓘ

ECO12 Berlin: Bernard Scherrer EDF Open Innovation (19 minutes) ⓘ

ENGL 3322 LECTURE 11 (56 minutes) ⓘ

Part 2 - The David and Lyn Silfen University Forum, 2012 (16 minutes) ⓘ

Denkfest: Pro Homeopathy? (36 minutes) ⓘ

Speaker: Edzard Ernst

What Can 'Friends of Syria' Do to Help Halt Killings? (10 minutes)

Dan Buettner: How to live to be 100+ (22 minutes) ⓘ

Given On: February 06, 2010

1

2

3

4

5

6

7

8

9

10

▶

1 to 20 of 200,212 lectures (~130,000 hours)

Weekly Emails

FREE TALK LIST

Want to learn something new? Sign up for our Monday morning list - an eclectic combination of our 3-5 videos and articles.

Email:

Sign me up!

2

# Goals

- Using machine learning on text
- Practical examples of Word2Vec in Scala
- Show uses of CUDA

# Agenda

- Proof of Concept: Email alerts
- Concept Search
- CUDA

## Papers

An empirical study of semantic similarity in WordNet and Word2Vec

<http://scholarworks.uno.edu/cgi/viewcontent.cgi?article=3003&context=td>

A Dual Embedding Space Model for Document Ranking

<https://arxiv.org/pdf/1602.01137v1.pdf>

## FindLectures.com Topic Alert

Keywords: scala

The following talks and articles were selected just for you:

### Articles

[Roadmap towards non-experimental macros](#)

[Getting Into Other People's Code](#)

[Configure Jupyter Notebook on Raspberry PI 2 for remote access and scala kernel install – log-IT.ro](#)

[Finatra Tutorial: Building Scalable Services The Twitter Way](#)

[Play 2.6.6 with sbt 1.0 support](#)

[Contest: scala-lang.org frontpage code snippet](#)

[Reactive Streams in Scala meet a Game Engine: part 6 - modularising the game logic](#)

### Videos

[5 Bullets to Scala Adoption \(Tomer Gabel, Israel\) \(50 minutes\)](#)

[Polyglot Programming with Python: Python/Scala Interop \(37 minutes\)](#)

[Scala Collections: Why Not? \(46 minutes\)](#)

This alert was generated by [Gary Sieling](#) using [FindLectures.com](#).

What did you think of this email?



Peace,  
Gary

P.S. You can report issues with these emails [on Github!](#)

# Email Alerts

## Personalized Weekly Tech News

Email

gary@garysieling.com

What topics interest you?

python machine learning

Suggestions: applications bootstrap browser campesina database excel format setup software

What topics would you prefer to avoid?

Suggestions: applications bootstrap browser campesina database excel format setup software

Submit

### Sample Talk Alert Email:

Here are this week's recommendations!



Lambda Days 2015 - Evelina Gabasova - Understanding cancer behaviour with F# (50 minutes)



Skdata: Data sets and algorithm evaluation protocols in Python; SciPy 2013 Presentation (24 minutes)



Rise of Scalable Machine Learning at Yahoo (39 minutes)

Peace,  
Gary

What did you think of this email?



## Concept Search

- Writing, NOT Code
- Excludes “writing css”, “writing php”
- Implies "poetry", "fiction", “copyediting”



## Concept Search

- Recipes, Vegetarian Food
- NOT Dairy
- All three might include "vegan cooking"
- Implies no milk, cheese

# Requirements

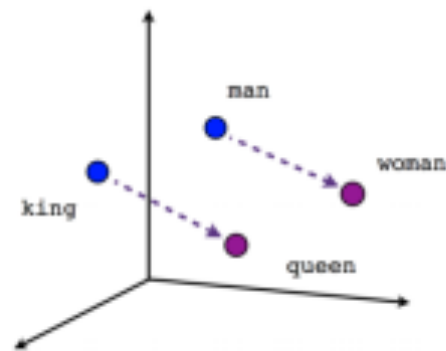
- Talks "about" the chosen topic
- Incorporate meaning – "Scala" + "Machine Learning" -> D|4j
  - May be a concept hierarchy
- Don't combine meaning if nothing in common (hiking, art)
- Don't send duplicate talks/articles (e.g. announcement from different publications)
- Choose a wide variety of talks (not 5 on type systems, etc)
- Bonus points for "negative" meanings (scala, but not monads)

## This is “search” problem

- Tokenize text
- Maybe mark known “entities”
- Filter / de-emphasize common terms / meanings
- Find the terms we should have searched for
- Search for those terms
- Re-rank / filter results

# Solution: Word2Vec

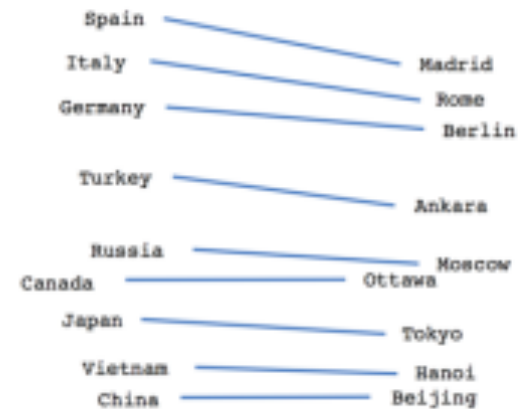
<https://github.com/idoio/wiki2vec>



Male-Female



Verb tense



Country-Capital

# Terms in context: Political Coding

<http://findlectures.com/?q=liberation>

## Faith in a Post-Modern World ►

**Description:** about political **liberation**? How can Christian faith be interpreted in a scientific and pluralist age?

## Henry Johnson comments on significance of the 1972 African Liberation Day March

**Description:** of the African **Liberation** Day March in Washington, 1972, and on-going significance of the march. Overall

## HIST 3322 LECTURE 6A (1.3 hours) Ⓞ

**Captions (00:08:31)** bit about the National **Liberation** Front

## Queer Studies Lecture Series - February 27, 2017 - Emily K. Hobson (48 minutes) Ⓞ

**Captions (00:03:07)** defined sexual **liberation** and radical

## GDL Primetime: Data Liberation (27 minutes) Ⓞ

**Captions (00:06:53)** Liberation and the Data **Liberation** Front all the way

## "Whale Wars" Ship a Terror Group? (2 minutes) Ⓞ

**Captions (00:00:18)** the Animal **Liberation** Front alf

## WORLD WAR II LIBERATION OF PARIS 1944 COMBAT DOCUMENTARY 78174 (13 minutes) Ⓞ

**Description:** in 1944. The **Liberation** of Paris (also known as the Battle for Paris) was a military conflict that took

## Paulo Freire & George Stoney Church Conversation (15 minutes) Ⓞ

**Description:** in 1996, Stoney and Freire discuss **Liberation** Theology, tolerance, large questions of life, and Freire's

# Terms in context: Context definitions

<http://findlectures.com/?q=quaker>

[This Separation Forced upon Us: Philadelphia's Free Quakers and the Culture of Revolution \(59 minutes\)](#)

**Description:** Despite their history of pacifism, Philadelphia **Quakers** were deeply entangled in the American

[Quakers Living Adventurously: The Library and Archives of the Society of Friends ►](#)

**Description:** Since the seventeenth century, members of the Religious Society of Friends - also known as **Quakers**

[Anthony Benezet, Father of Atlantic Abolitionism \(60 minutes\)](#)

**Description:** to convince his **Quaker** brethren that slave-owning was not consistent with Christian doctrine. Benezet and his

[The Life of Herbert Hoover: Fighting Quaker 1928-1933 \(52 minutes\) ⓘ](#)

**Captions (00:21:15)** economically than most of the Western world. The first **Quaker** president, the first born

[2016 Stephen G. Cary Lecture \(1.5 hours\) ⓘ](#)

**Captions (01:24:39)** which I understood **Quaker**

[David Holmes - Religion and Watergate - 03/17/15 \(60 minutes\) ⓘ](#)

**Captions (00:06:28)** was raised a **Quaker**

[Book TV at UCLA: Lane Hirabayashi, "A Principled Stand" \(10 minutes\)](#)

**Captions (00:06:14)** he was already a **Quaker** because he was

[President Reagan's Radio Address on Tax Reforms from Camp David, Maryland on June 1, 1985 \(5 minutes\) ⓘ](#)

**Captions (00:00:10)** the **Quaker** whose patience was sorely

[Philip Gulley: 2015 Stephen G. Cary Memorial Lecture \(1.4 hours\) ⓘ](#)

**Description:** **Quaker** Pastor speaks about programmed and un-programmed **Quaker** meetings.

# Training Vectors

Was raised a Quaker

["was", "raised", "a", "religious", "since", "the", "whose", "patience"]

[1, 1, 1, 0, 0, 0, 0, 0]

The Quaker whose patience was

["was", "raised", "a", "religious", "since", "the", "whose", "patience"]

[1, 0, 0, 0, 0, 1, 1, 1]

## Word2Vec Output

$P(\text{Term} \mid \text{context})$

Or

$P(\text{Context} \mid \text{Term})$



## Example: Vector Addition

Gloria Steinem - Person + Ideology  $\sim$  =

1. Marxist Feminism
2. Radical Feminism
3. Feminist Movement
4. Feminist Theory

# Suggested Search

## Searches related to Pennsylvania:

<a href="#">Pennsylvania</a>	<a href="#">Philadelphia, Pennsylvania</a>
<a href="#">Montgomery County, Pennsylvania</a>	<a href="#">York County, Pennsylvania</a>
<a href="#">Bucks County, Pennsylvania</a>	<a href="#">Philadelphia</a>
<a href="#">Pittsburgh, Pennsylvania</a>	<a href="#">Harrisburg, Pennsylvania</a>
<a href="#">Berks County, Pennsylvania</a>	<a href="#">Lancaster County, Pennsylvania</a>

## Searches related to Bayard Rustin:

<a href="#">Rustin</a>	<a href="#">A. Philip Randolph</a>
<a href="#">Bayard</a>	<a href="#">Tom Kahn</a>
<a href="#">civil-rights</a>	<a href="#">Coretta Scott King</a>
<a href="#">Ella Baker</a>	<a href="#">A. J. Muste</a>
<a href="#">Roy Wilkins</a>	<a href="#">James L. Farmer, Jr.</a>

## Searches related to Stokely Carmichael:

<a href="#">Stokely</a>	<a href="#">SNCC</a>
<a href="#">Student Nonviolent Coordinating Committee</a>	<a href="#">James Bevel</a>
<a href="#">Malcolm X</a>	<a href="#">Diane Nash</a>
<a href="#">(SNCC)</a>	<a href="#">Fannie Lou Hamer</a>
<a href="#">Robert Parris Moses</a>	<a href="#">Black Power</a>

## Searches related to Social Justice:

<a href="#">Social change</a>	<a href="#">Social responsibility</a>
<a href="#">Social equality</a>	<a href="#">Teaching for social justice</a>
<a href="#">Social</a>	<a href="#">Social movement</a>
<a href="#">Social exclusion</a>	<a href="#">Social market economy</a>
<a href="#">social justice</a>	<a href="#">justice</a>

## Searches related to Philadelphia Eagles:

<a href="#">Chicago Bears</a>	<a href="#">New York Giants</a>
<a href="#">Carolina Panthers</a>	<a href="#">Detroit Lions</a>
<a href="#">New York Jets</a>	<a href="#">Arizona Cardinals</a>
<a href="#">New England Patriots</a>	<a href="#">National Football League</a>
<a href="#">Atlanta Falcons</a>	<a href="#">Cincinnati Bengals</a>

## Searches related to Machine Learning:

<a href="#">Machine Learning</a>	<a href="#">Pattern recognition</a>
<a href="#">Neural network</a>	<a href="#">Support vector machine</a>
<a href="#">Artificial neural network</a>	<a href="#">Knowledge representation</a>
<a href="#">machine learning</a>	<a href="#">Evolutionary computation</a>
<a href="#">Subsumption</a>	<a href="#">Multi-agent system</a>

## Example: Data Format

```
{
  "word": "zułus"
  "count": 30,
  "syn0": [
    -0.064, 0.118, 0.031, 0.163, 0.019, 0.197, 0.097, -0.139, -0.055, 0.155,
    -0.033, -0.252, -0.029, 0.119, 0.007, -0.017, 0.187, 0.017, 0.058, -0.097,
    -0.255, -0.159, -0.053, -0.090, -0.118, 0.119, 0.068, 0.025, 0.160, -0.035,
    -0.216, 0.065, 0.017, 0.038, -0.068, 0.101, 0.090, 0.089, -0.023, 0.265,
    -0.161, -0.178, -0.362, 0.016, 0.226, -0.070, -0.079, 0.040, 0.368, -0.150
  ],
  "syn1": [
    0.312, 0.379, 0.168, -0.371, -0.094, 0.218, -0.022, -0.051, 0.003, -0.010,
    0.233, -0.005, -0.037, 0.105, 0.025, -0.040, -0.127, .201, 0.175, 0.277,
    0.185, -0.219, -0.504, -0.187, 0.069, 0.041, 0.237, -0.245, 0.067,
    -0.186, 0.127, 0.235, -0.262, -0.020, -0.152, 0.007, -0.346, 0.008, -0.173,
    -0.267, -0.049, 0.051, 0.087, 0.046, -0.059, 0.147, 0.024, 0.032, -0.403,
    0.019
  ]
}
```

## Example: Similarity

Number from  $[0, 1]$

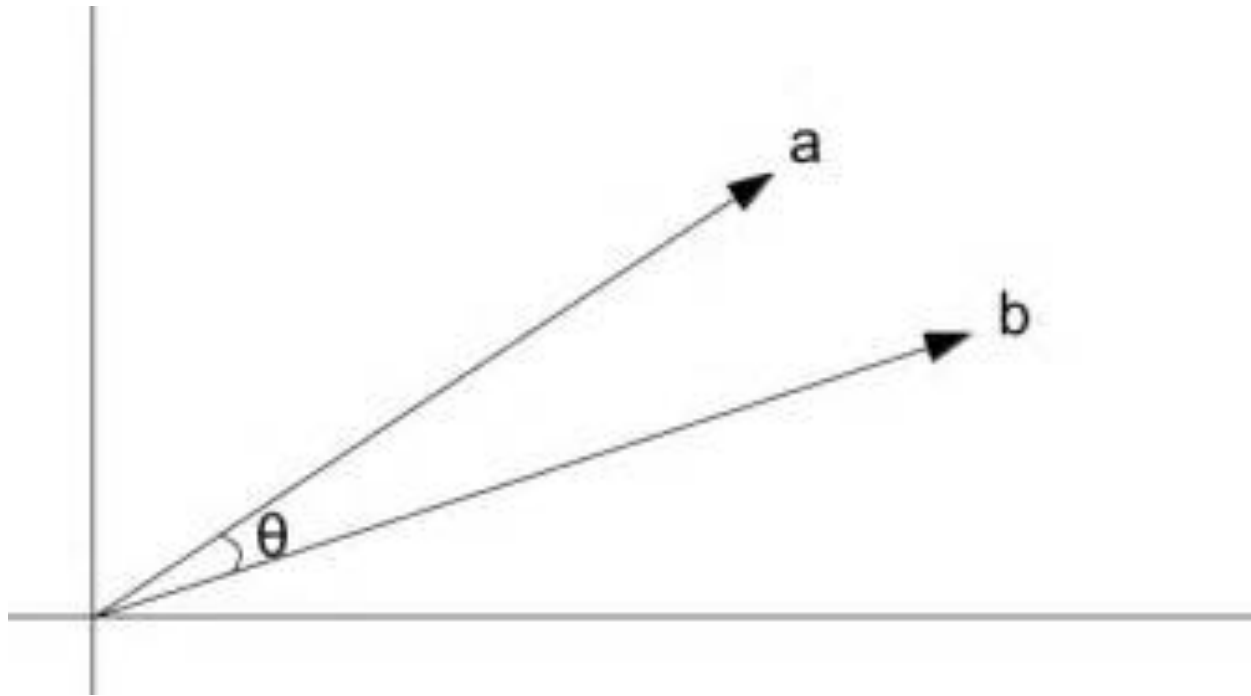


Image credit: <https://engineering.aweber.com/cosine-similarity/>

## Operation 1: “Similarity”

```
def cosineSimilarity(  
    a: INDArray,  
    b: INDArray  
): Double = {  
    Transforms.cosineSim(a, b)  
}
```

## INDArray

- Similar to numpy array
- Implementation depends on dependency:

libraryDependencies +=

"org.nd4j" % "nd4j-cuda-8.0-platform" % nd4jVersion

libraryDependencies +=

"org.nd4j" % "nd4j-native" % nd4jVersion

# CUDA

- Specialized instruction set in video cards / GPUs
- Requires NVIDIA SDK and a recent card (\$100-\$xx,xxx)
- Available on AWS
- Deeplearning4j: JVM libraries for machine learning
- Nd4j/nd4s: matrix algebra on large arrays

# CUDA: example C code

```
__global__ void coalescedMultiply(float *a, float *c, int M)
{
    __shared__ float aTile[TILE_DIM][TILE_DIM],
    transposedTile[TILE_DIM][TILE_DIM];
    int row = blockIdx.y * blockDim.y + threadIdx.y;
    int col = blockIdx.x * blockDim.x + threadIdx.x;
    float sum = 0.0f;
    aTile[threadIdx.y][threadIdx.x] = a[row*TILE_DIM+threadIdx.x];
    transposedTile[threadIdx.x][threadIdx.y] =
        a[(blockIdx.x*blockDim.x + threadIdx.y)*TILE_DIM +
        threadIdx.x];
    __syncthreads();
    for (int i = 0; i < TILE_DIM; i++)
        sum += aTile[threadIdx.y][i] * transposedTile[i][threadIdx.x];
    c[row*M+col] = sum;
}
```



## Ways to obtain GPUS

- Buying
- Renting
  - AWS (\$0.90/hr)

Name	GPUs	vCPUs	RAM (GiB)	Network Bandwidth	Price/Hour*	RI Price / Hour**
p2.xlarge	1	4	61	High	\$0.900	\$0.425
p2.8xlarge	8	32	488	10 Gbps	\$7.200	\$3.400
p2.16xlarge	16	64	732	20 Gbps	\$14.400	\$6.800

# Training Word2Vec

```
val vec =  
    new Word2Vec.Builder()  
        .minWordFrequency(5)  
        .iterations(1)  
        .layerSize(100)  
        .seed(42)  
        .windowSize(5)  
        .iterate(sentenceliterator)  
        .tokenizerFactory(tokenizer)  
        .build  
vec.fit();
```

# How do you tell if your code is running - GPU



# How does this affect word2vec

- Dl4j Demo project: 72 minutes (CPU)
- Dl4j Demo project: 41 minutes (GPU)

# Most Similar

....

Definining ops we can use – should this be sooner?

## Operation 2: Compute a document mean

```
def getWordVectorsMean(tokens: List[String]): INDArray = {  
    val words = tokens.filter(  
        model.getWordVector(_) != null  
    ).sorted  
  
    model.getWordVectorsMean(  
        words.asJavaCollection  
    )  
}
```

## Nd4s / Nd4j

- Everything is one long array, with dimensions (like numpy)
- Create one with a big iterator
- Easy to reshape
- Parallelism – min 32 cores, all following same path

# Problem: Suggestions

By the next search?

## Searches related to ron chernow

ron chernow **website**

ron chernow **contact**

ron chernow **grant**

ron chernow **alexander hamilton**

ron chernow **titan**

ron chernow **washington**

ron chernow **awards**

ron chernow **alexander hamilton pdf**



# Problem: Noise

## Personal background [\[ edit \]](#)

---

Ron Chernow has received honorary degrees from [Long Island University](#), [Marymount Manhattan College](#), [Hamilton College](#), [Washington College](#), and [Skidmore College](#).<sup>[3]</sup>

# Nd4s – Make an array

```
val data: Seq[Double] =  
  Seq(  
    words.flatMap(  
      (w) => wordVectors(w)  
    ),  
    words.flatMap(  
      (w) => Seq.iterate(1, widthOfWordVector)((idx: Int) => termFrequencies(w)).map(  
        (vv: Int) => vv.toDouble  
      )  
    ),  
    words.flatMap(  
      (w) => Seq.iterate(1, widthOfWordVector)((idx: Int) => documentFrequencies(w)).map(  
        (vv: Int) => vv.toDouble  
      )  
    )  
  ).flatten
```

# Nd4s – Computation of TF\*IDF average

```
val modeVectors = arr.reshape(modes, widthOfWordVector * numWords)
    val scores = modeVectors(0 -> 1)
    val tf = modeVectors(1 -> 2)
    val df = modeVectors(2 -> 3)

    val weighted = scores * tf / df

    val wordVects = weighted.reshape(numWords, widthOfWordVector)
    // this is the weighted average

    wordVects.sum(0) / numWords
```

```
// TODO is this any better?
```

## "Synonym" Discovery Example

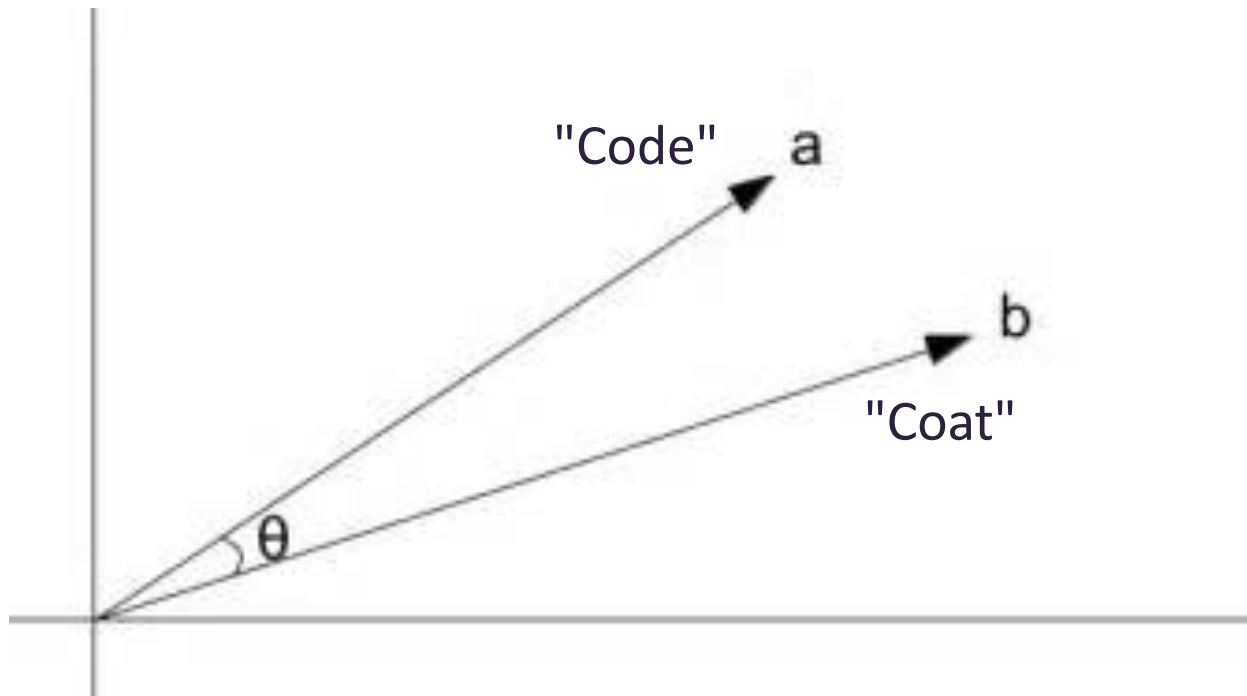


Image credit: <https://engineering.aweber.com/cosine-similarity/>

## Word2Vec – Build a Full Text Query

```
List("python", "machine", "learning").map(
  (queryTerm) =>
    "(" +
      model.wordsNearest(
        List(queryTerm), // positive terms
        List(), // negative terms
        25
      ).map(
        (nearWord) =>
          "transcript:" + term2 +
            "^" + model.similarity(nearWord, term2)
      ).mkString(" OR ")
    + ")"
  ).mkString(" AND ")
```

## Visual – Nearest terms

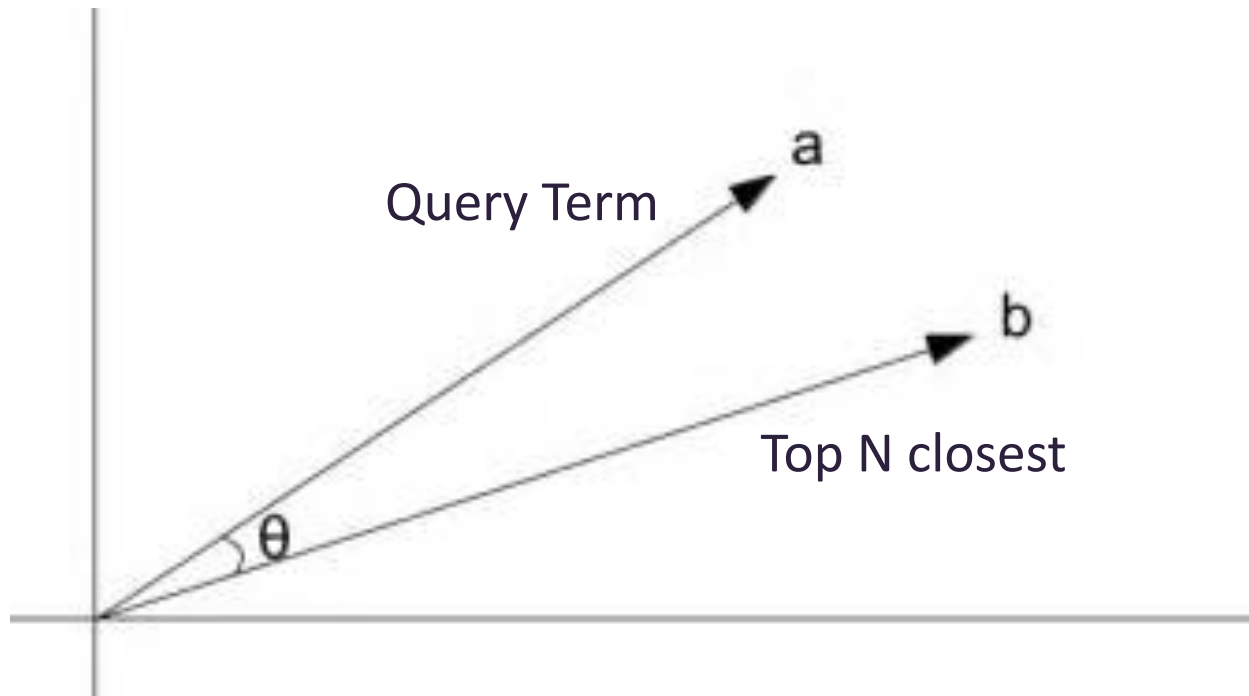


Image credit: <https://engineering.aweber.com/cosine-similarity/>

## Example – Query (“Python + Machine Learning”)

```
title_s:python^10 OR title_s:"machine learning"^10 ...  
(title_s: software^1.21 OR title_s:database^1.20 OR title_s:format^1.18  
title_s:applications^1.14 OR title_s:browser^1.14 OR title_s:setup^1.13  
title_s:bootstrap^1.13 OR title_s:in-class^1.13 OR title_s:campesina^1.12 OR  
title_s:excel^1.12 OR title_s:hardware^1.11 OR title_s:programming^1.11 OR  
title_s:api^1.11 OR title_s:prototype^1.11 OR title_s:middleware^1.11 OR  
title_s:openstreetmap^1.10 OR title_s:product^1.10 OR title_s:app^1.09 OR  
title_s:hbp^1.09 OR title_s:programmers^1.09 OR title_s:application^1.09 OR  
title_s:databases^1.09 OR title_s:idiomatic^1.09 OR title_s:spreadsheet^1.09  
OR title_s:java^1.09 ...  
AND (...)
```

## Results (Python + Machine Learning + BM25)

Python for Data Analysis

How To Get Started With Machine Learning? | Two Minute Papers

The /r/playrust Classifier: Real World Rust Data Science

Andreas Mueller - Commodity Machine Learning

A Gentle Introduction To Machine Learning

A full Machine learning pipeline in Scikit-learn vs in scala-Spark

Hello World - Machine Learning Recipes #1

Visual diagnostics for more informed machine learning

Lab to Factory: Robust Machine Learning Systems

**Machine Learning with Scala on Spark by Jose Quesada**



## Word2Vec – “Writing”

Issues Related to the Teaching of Creative Writing  
Is Nonfiction Literature?

"Oh, you liar, you storyteller": On Fibbing, Fact and Fabulation  
The Value of the Essay in the 21st Century

**Re writing Re reading Re thinking – Web Design in Words**

Aspen New York Book Series: The Art of the Memoir

Cheryl Strayed: "Wild"

Siri Hustvedt in Conversation with Paul Auster

Mary Karr: The 2016 Diana and Simon Raab Writer-in-Residence  
History, Memory, and the Novel

# Aboutness

Re-sorting top 100 documents

```
val queryMean = model.getWordVectorsMean(List("writing"))  
val mean = model.getWordVectorsMean(NLP.getWords(document._1))  
val distance = Transforms.cosineSim(vec._2, queryMean)
```

5 min 45 seconds @ 16 parallel threads

## Visual – Aboutness

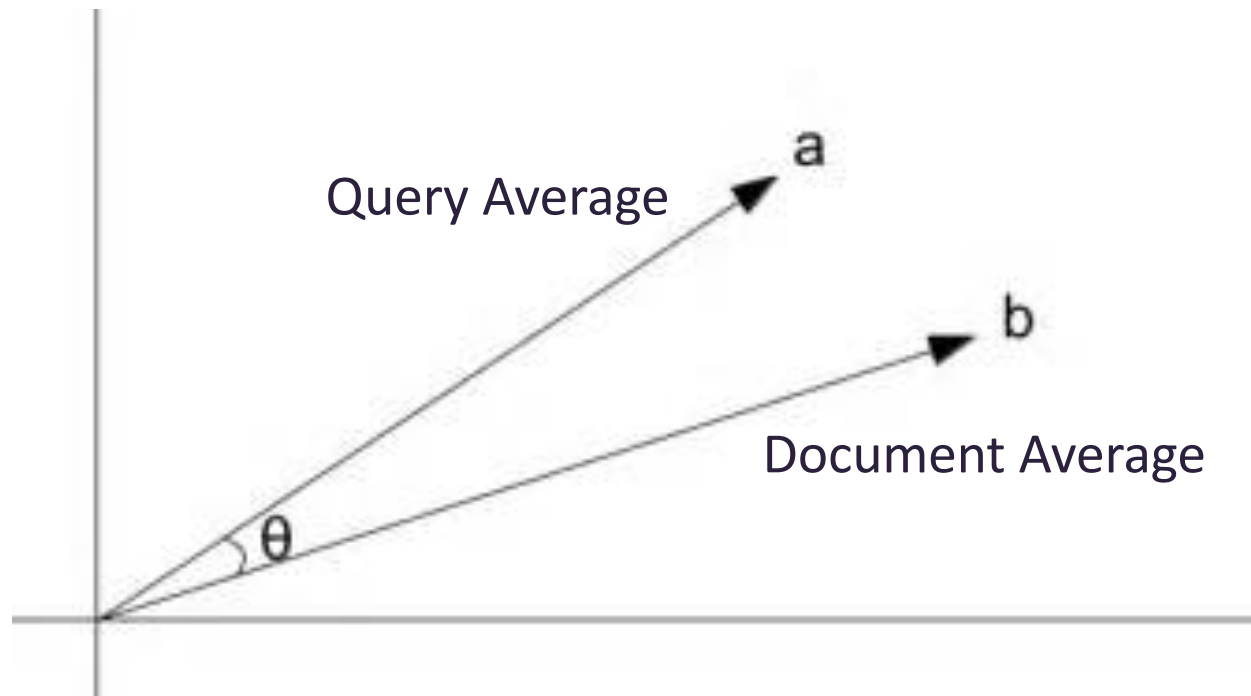


Image credit: <https://engineering.aweber.com/cosine-similarity/>

## Aboutness - Results

Issues Related to the Teaching of Creative Writing: 0.43

Autobiography: 0.41

Contemporary Indian Writers: The Search for Creativity: 0.41

Marjorie Welish: Lecture: 0.40

History and Literature: The State of Play: A Roundtable Discussion: 0.40

Critical Reading of Great Writers: Albert Camus: 0.40

Daniel Schwarz: In Defense of Reading: 0.39

The Journey To The West by Professor Anthony C. Yu: 0.39

Blogs, Twitter, the Kindle: The Future of Reading: 0.39

# Word2Vec + Overlapping Search Terms

Python, Programming vs Art, Hiking

```
terms.map(  
  (term1) =>  
    terms.map(  
      (term2) => (term1, term2)  
    )  
).flatten.filter(  
  (tuple) => tuple._1 < tuple._2  
).map(  
  (tuple) =>  
    (tuple._1, tuple._2, w2v.model.get.similarity(tuple._1, tuple._2))  
)
```

## Visual – Overlapping Search Terms

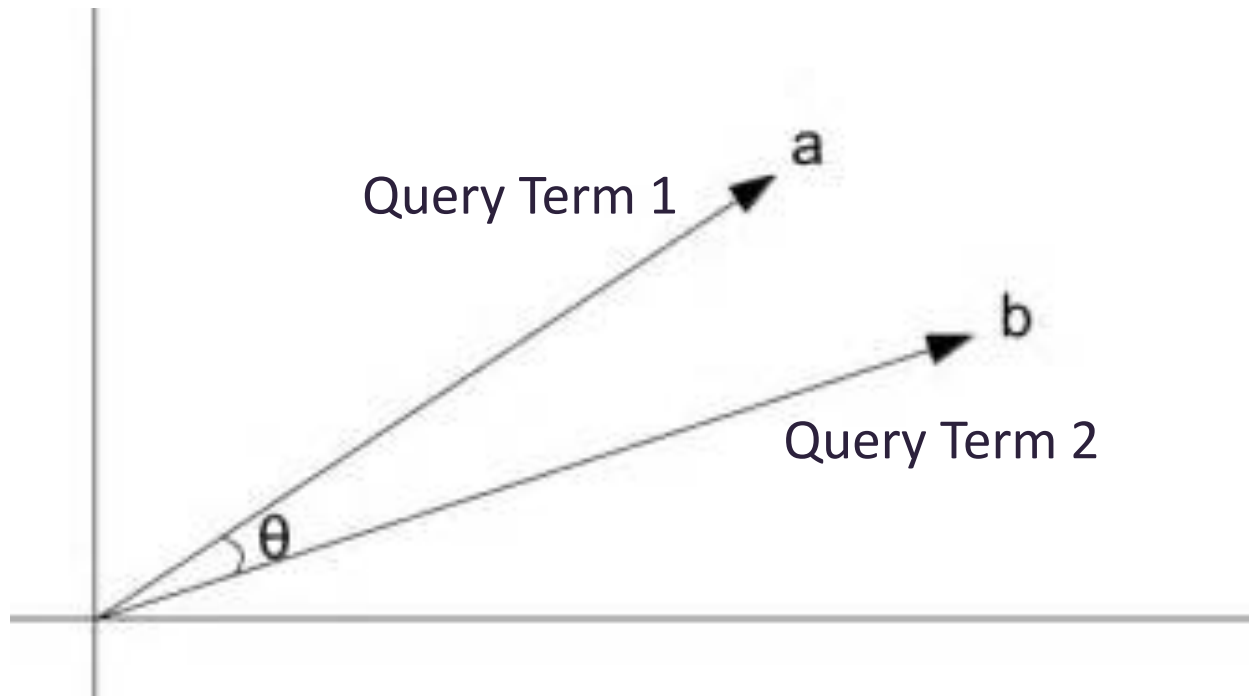


Image credit: <https://engineering.aweber.com/cosine-similarity/>

## Word2Vec + Overlapping Search Terms

Python, Programming

programming $\leftrightarrow$ python: 0.61

(python AND programming)

Hiking, Art

art $\leftrightarrow$ hiking: 0.10

(hiking OR art)

## Topic Diversity

### Writing

[A Conversation with David Gerrold, Writer of Star Trek: The Trouble with Tribbles - Teletalk \(58 minutes\)](#)

[Star Trek: Science Fiction to Science Fact - STEM in 30 \(28 minutes\)](#)

### Python

[Pythons Positive Press Pumps Pandas](#)

[Why is Python Growing So Quickly? - Stack Overflow Blog](#)

[Python explosion blamed on pandas](#)



## Visual – Topic Diversity

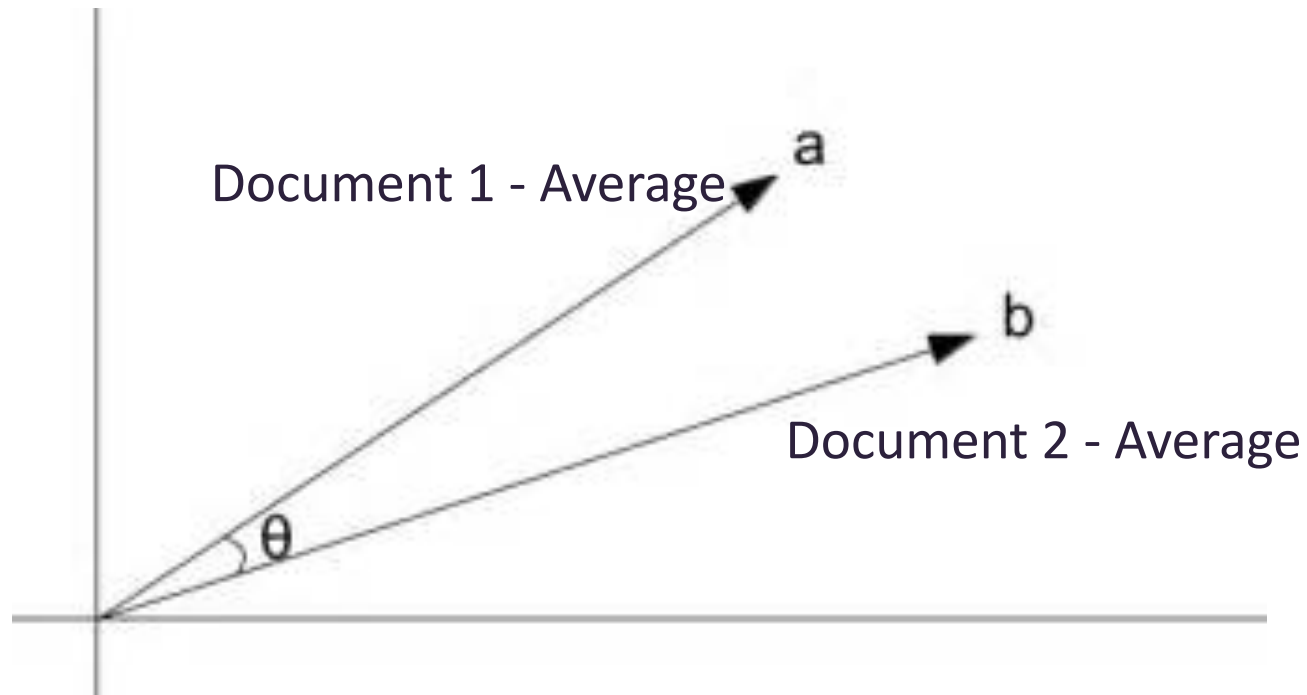


Image credit: <https://engineering.aweber.com/cosine-similarity/>

Pick one, find the least related (Python + Pandas)

**Python explosion blamed on pandas: 1.0**

Considering Python's Target Audience: 0.97

Animated routes with QGIS and Python: 0.97

I can't get some SQL to commit reading data from a database: 0.97

Using Python to build an AI Twitter bot people trust: 0.96

Getting a Job as a Self-Taught Python Developer: 0.96

Download and Process DEMs in Python: 0.96

How to mine newsfeed data and extract interactive insights in Python: 0.94

Differential Equation Solver In MATLAB, R, Julia, Python, C, Mathematica, Maple, and Fortran: 0.86

My personal data science toolbox written in Python: 0.75

1 min 30 seconds @ 16 parallel threads

## Technique - Summary

- Get top X results, re-shuffle
- More computing resources + data -> higher relevance

## Where Word2Vec Works

- Synonym generation
- Improve recall
- Search suggestions
- Incorporate secondary dataset (e.g. for enterprise search, privacy)

# Why Scala?

- Ecosystem: Lucene, Spark
- Dependency Management

# Performance

- Models take 1-2 weeks to train
- Some of computations take minutes, which would not work in a search engine
- Changes:
  - Pre-compute tokens (e.g. use Lucene)
  - Pre-compute averages (don't naturally store in Lucene)
  - Hazelcast

# How do you tell if your code is running on a GPU (Spark + Deeplearning4j)

- 15:17:27,828 INFO ~ Loaded [CpuBackend] backend
- 15:17:28,008 INFO ~ Number of threads used for NativeOps: 4
- 15:17:29,182 INFO ~ Number of threads used for BLAS: 4
- 15:17:29,185 INFO ~ Backend used: [CPU]; OS: [Windows 10]
- 15:17:29,185 INFO ~ Cores: [8]; Memory: [3.6GB];
- 15:17:29,185 INFO ~ Blas vendor: [MKL]
- 15:17:34,546 INFO ~ Using Spark Local

# CUDA

- Switch between CPU and GPU by changing sbt configuration:
- Threading resources. Execution pipelines on host systems can support a limited number of concurrent threads. Servers that have four hex-core processors today can run only 24 threads concurrently (or 48 if the CPUs support HyperThreading.) By comparison, the smallest executable unit of parallelism on a CUDA device comprises 32 threads (termed a warp of threads). Modern NVIDIA GPUs can support up to 1536 active threads concurrently per multiprocessor (see Section F.1 of the CUDA C Programming Guide). On GPUs with 16 multiprocessors, this leads to more than 24,000 concurrently active threads.



## Hazelcast

- Just videos – 241.8 minutes
- Nothing cached, but hazelcast- 76 minutes
- On query combos – 234 minutes
- Adding Hazelcast on queries - 62.091
- After all cached – 2.38
- Move word2vec model from spinner to SSD:

# jCuda

```
def memory = {  
  cuInit(0)  
  val device = new CUdevice  
  JCudaDriver.cuDeviceGet(device, deviceId)  
  
  val total = Array(0L)  
  val free = Array(0L)  
  cuInit(0)  
  cuDeviceGet(device, deviceId)  
  
  val context = new CUcontext  
  cuCtxCreate(context, 0, device)  
  cuMemGetInfo(free, total)  
  
  cuCtxDestroy(context)  
  
  (total(0), free(0))  
}
```

# Tokenize - Lucene

```
def getTokens(text: String): List[String] = {  
  val result = new util.ArrayList[String]()  
  val analyzer: Analyzer = new StandardAnalyzer()  
  
  val stream: TokenStream = analyzer.tokenStream(null, new StringReader(text))  
  stream.reset()  
  
  while (stream.incrementToken) {  
    result.add(stream.getAttribute(classOf[CharTermAttribute]).toString())  
  }  
  
  import scala.collection.JavaConversions._  
  result.toList  
}
```

# Other Lessons

- Inventing your own math does not work
  - High-dimensional “objects” do not follow your intuition like 2D/3D
  - Floating point math not associative
- Math in papers is untyped
  - “Distance” between two vectors – cosine, euclidean, manhattan?
  - vs. Probability curves
  - Unlike Physics ( types naturally compose,  $\text{kg} \cdot \text{m}^2 \cdot \text{s}^{-2}$  )
- Follow a paper
  - Nearly impossible to test on your own
  - Almost no one publishes code

# Next Idea...

React App

localhost:3001/topic/social%20justice


Search...

Topic > Social Justice


Related Topics: [Socialism](#), [Feminism](#), [Black History](#), [Liberation Theology](#), [Non-profit management](#)

Social justice is a concept of fair and just relations between the individual and society. The lectures below include major themes, influential speakers, and primary source materials.


Historic Speeches




Barack Obama at NAACP Annual Conference




CAIR Joins U.S. Muslim Leaders in Condemning Boko Haram



Kennedy Endorsements of Barack Obama




Morgan State 2012 Commencement Address: Dr. Shirley Ann Jackson




Perspectivas Latinas - Presencia Michoacana in the Midwest 2014


Economic Justice




Video: Full CAIR D.C. News Conference Calling on Ben Carson to Withdraw from Presidential Race




Book TV: Gordon Mantler, "Power to the Poor"



Virginia Eubanks: Deconstructing the Digital Divide




Killer Mike on Russell Simmons' "Bernie Insensitive to Plight of Black People"




Reproductive cloning  
Marcy Darnovsky: Should We Genetically Modify Our Children?


Civil Rights Movement




For Jobs and Freedom: 50 Years and




BookTV: 2011 Brooklyn Book Festival -



For Jobs and Freedom: A Black Nouveau Special | Program



Adam Green: The Process of Diversity



Jane Wales and Larry Brilliant: Welcome and Conference...

localhost:3001/view/198695

# CUDA Surprises

- High end GPUs don't do video
- A ton of people are using these for bitcoin mining (see local craigslist)
- CUDA uses a lot of CPU
- Floating-Point Math Is Not Associative
- "...the peak theoretical memory bandwidth of the NVIDIA Tesla M2090 is 177.6 GB/sec:  $(1.85 \times 10^9 \times (384/8) \times 2) / 10^9 = 177.6 \text{ GB/sec}$ "
- ".... the peak theoretical bandwidth between host memory and device memory (8 GB/s on the PCIe x16 Gen2).
- "...if, switch, do, for, while significantly affect throughput ... The different execution paths must be serialized, since all threads of a warp share a program counter; this increases the total number of instructions executed for this warp"

## Resources

- "Relevant Search"
- "Deep Learning – A Practitioner's Approach"
- Deeplearning4j
- Gensim
- <https://github.com/DiceTechJobs/ConceptualSearch>
- [https://www.reddit.com/r/datasets/comments/3mg812/full\\_r\\_eddit\\_submission\\_corpus\\_now\\_available\\_2006/](https://www.reddit.com/r/datasets/comments/3mg812/full_r_eddit_submission_corpus_now_available_2006/)

[FindLectures.com](http://findlectures.com)

Weekly Emails with Lunch and Learn Suggestions

<http://findlectures.com/emails>



Next installment:

Java Users Group In February 2018

**“GPU Programming for Java Developers”**

Contact:

@garysieling

@findlectures

[gary@garysieling.com](mailto:gary@garysieling.com)

<https://www.findlectures.com>

<https://www.garysieling.com>

<https://github.com/garysieling/>