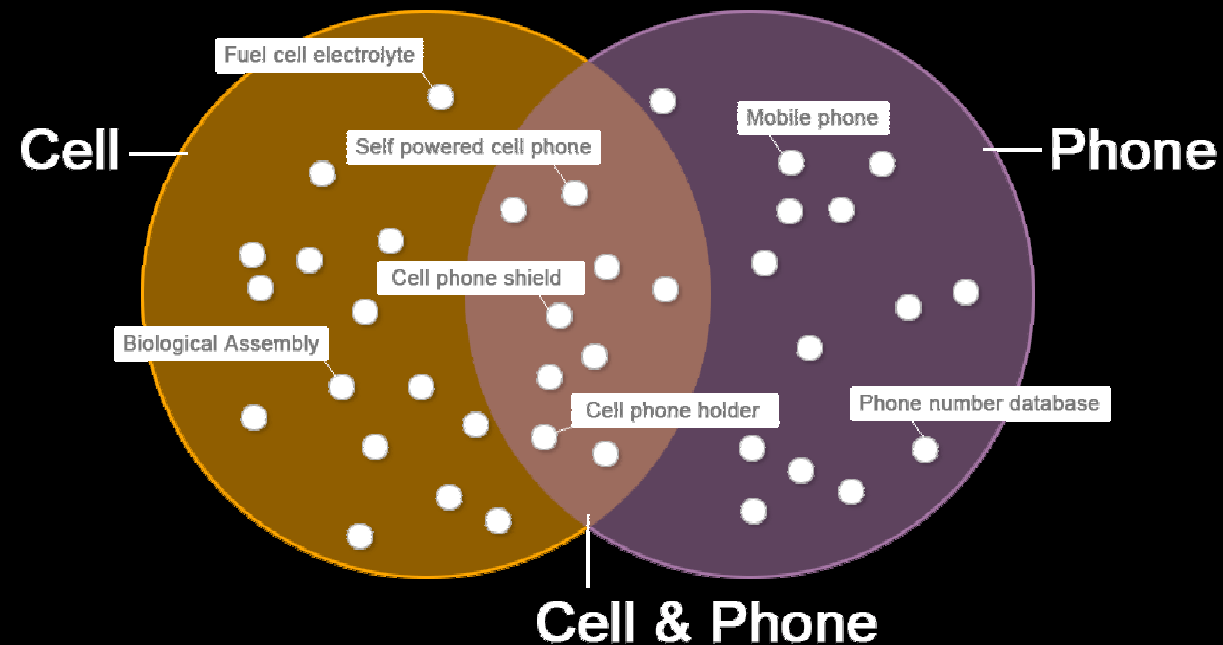# Using Conceptual Search in Scientific, Financial and Intellectual Property Databases

James Ryley, Ph.D.

# Conventional (Boolean) Search

- Works on the basis of set intersection. There is no "understanding" of document language.

# Recall and Precision

Recall and precision suffer when language is not understood.

Consider that, for the search query "cell AND phone":

Documents containing the terms "mobile phone" or "hand phone" (common in Asia) will be missed.

# Recall and Precision

Switching to "cell OR phone" leads to a result set containing documents about:

Terrorist cells
Jail cells
Electrolytic cells
Sony Playstation Cell™ MicroProcessor
Stem Cells

and more…

# Improving Recall and Precision in Boolean Search

- The solution is "simple," (though creates ranking nightmares) just create a query like this:

  (phone OR telephone OR radio) AND (mobile OR cell OR cellular OR hand OR wireless OR radio OR satellite OR portable OR … plus the terms for your specific interest in cell phones

# Latent Semantic Analysis (LSA)

- Searches conceptually, not literally, enhancing recall and precision

- Uses a spatial representation of documents to enhance precision through accurate ranking

# LSA Concept Mapping

- Statistical techniques are used to relate terms together.
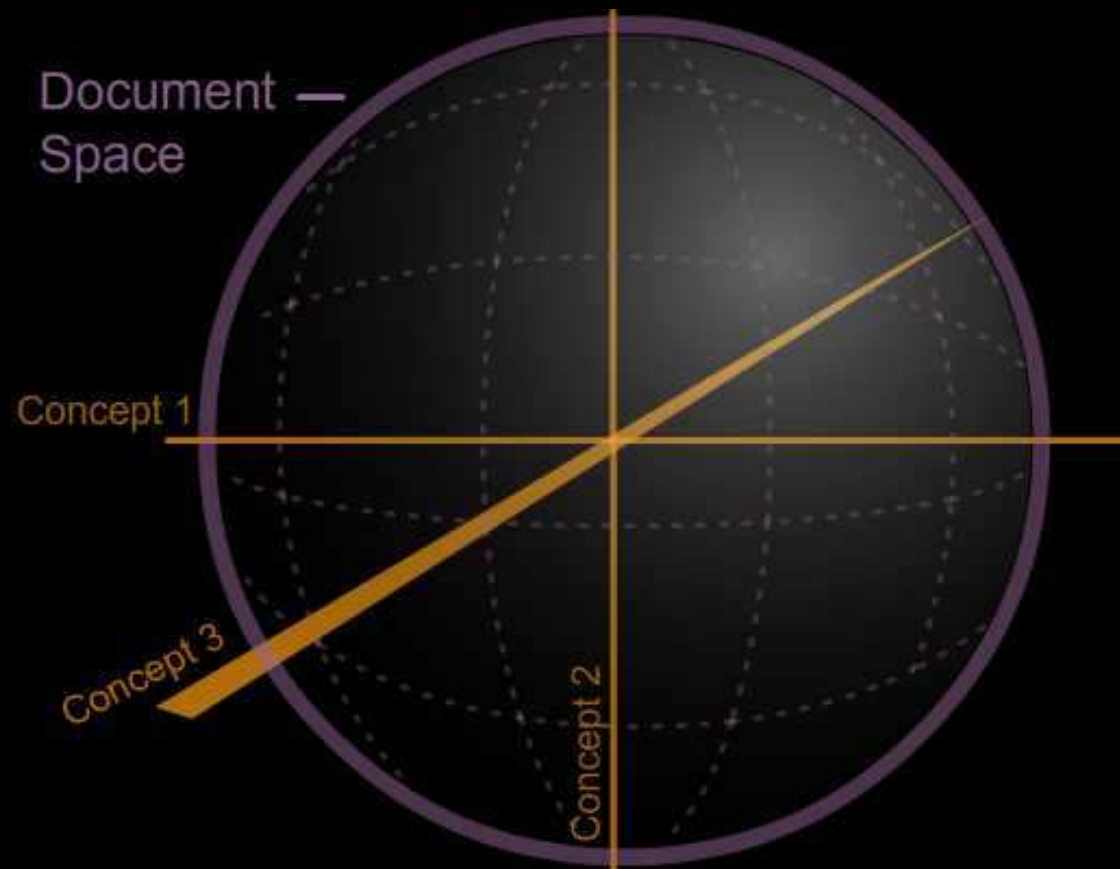- The result is the ability to find documents that do not contain the literal search terms.

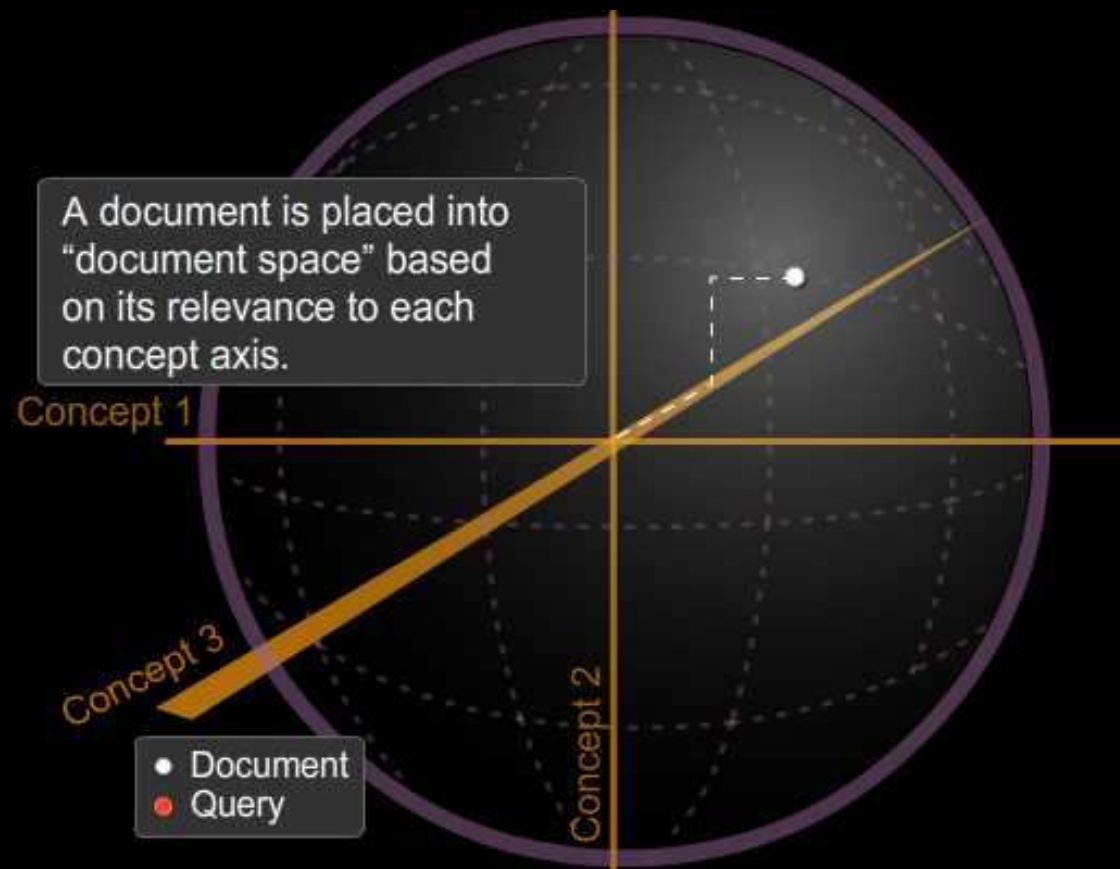| Terms |
| --- |
| cell, cellular, hand, mobile, phone, telephone, radio |

→

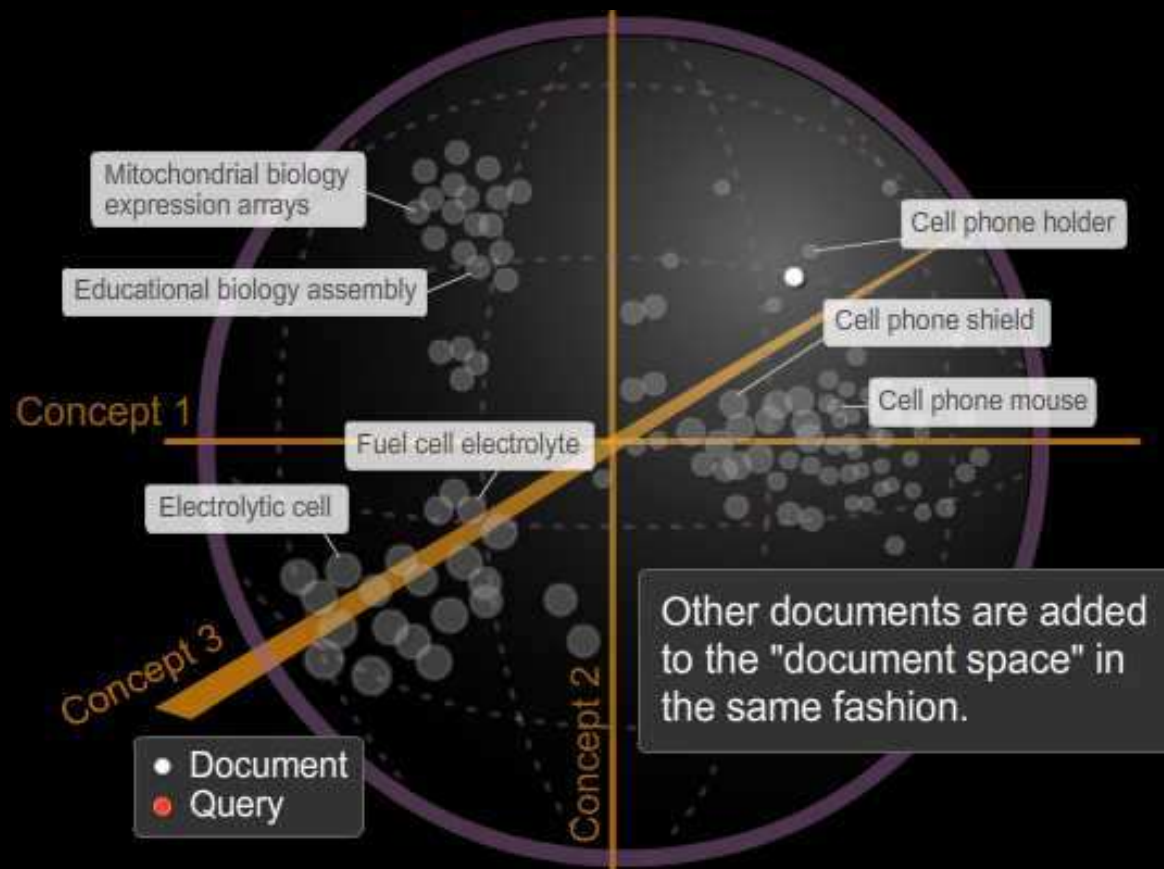| Concept |
| --- |
| "cellular phone" |

# The Spatial Document Model

# Filling Up Document Space

# Spatial Clustering

# Ranking Using Angles

Cell (Biology)

Cell (Phones)

Concept 1

Cell (Electrolytic)

Concept 3

Concept 2

To find relevant documents, the angle between the query and each document is compared.

- Document
- Query

# LSA Summary

- The system "understands" concepts, not just words

- Polysemous words ("cell" in "cell phone" vs. "stem cell"), are appropriately ignored

- Results are ordered accurately by relevance

**LSA has high recall, and high precision.**

# If LSA is so Great, Why Isn't Everyone Using It?

- Mathematical and logical complexity
- Large computational demands
- *Lack of transparency:  Can we get by this??*

**SUMOBRAIN**
Big, powerful, smart searching

Next-Generation Patent Searching

Thanks!

Contact info:

james.ryley@sumobrain.com

www.SumoBrain.com