

Comparison of Hyper-dimensional LSA Spaces for Semantic Differences

John C. Martin

Dissertation Defense

20 May 2016

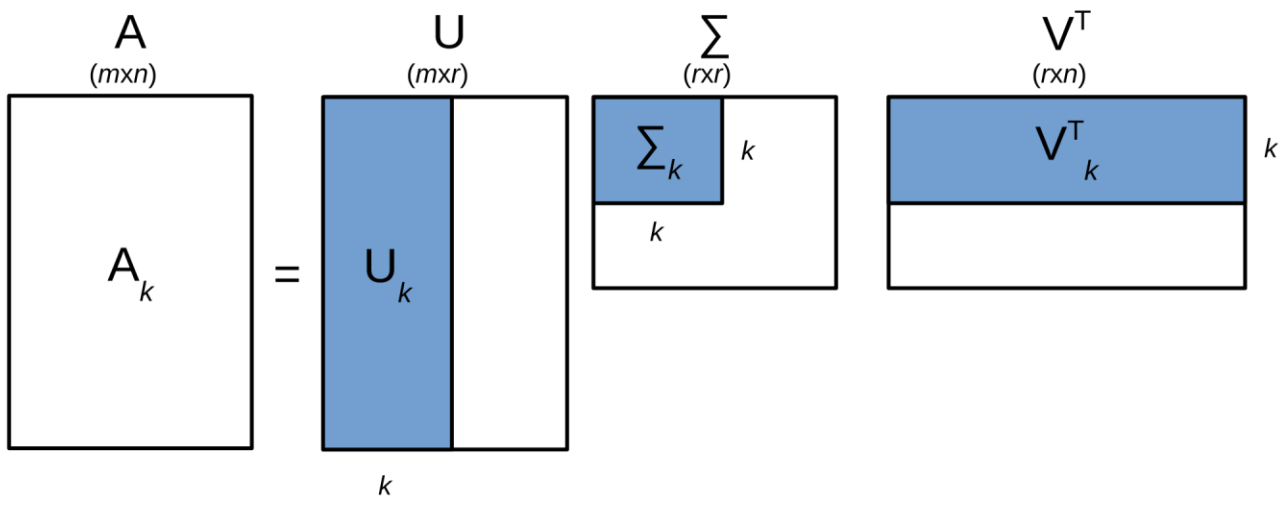
Overview

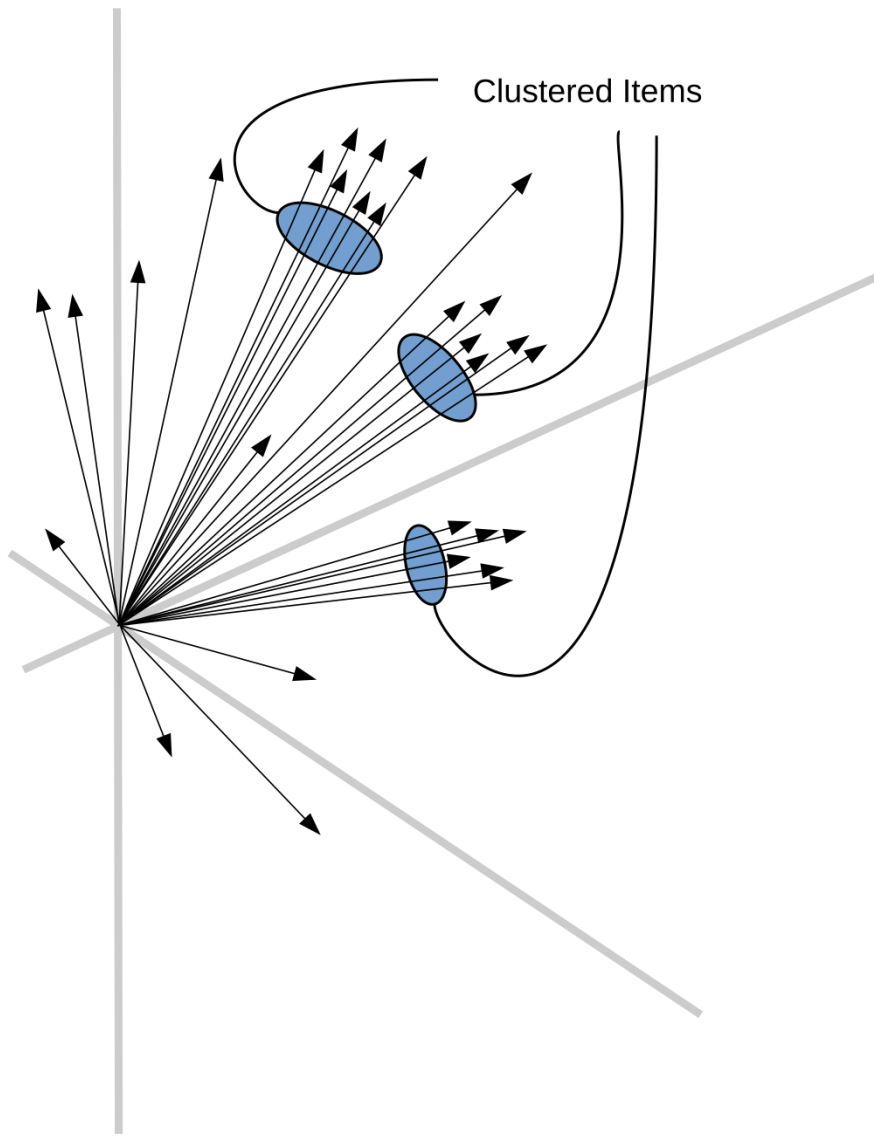
- Review LSA model of learning
 - What is meaning?
- Measures
- Experiments
- Semantic Measurement Model
- Q & A

The LSA Model of Learning

- Orthogonal Axes
- Dimensionality Reduction
- Mapping System

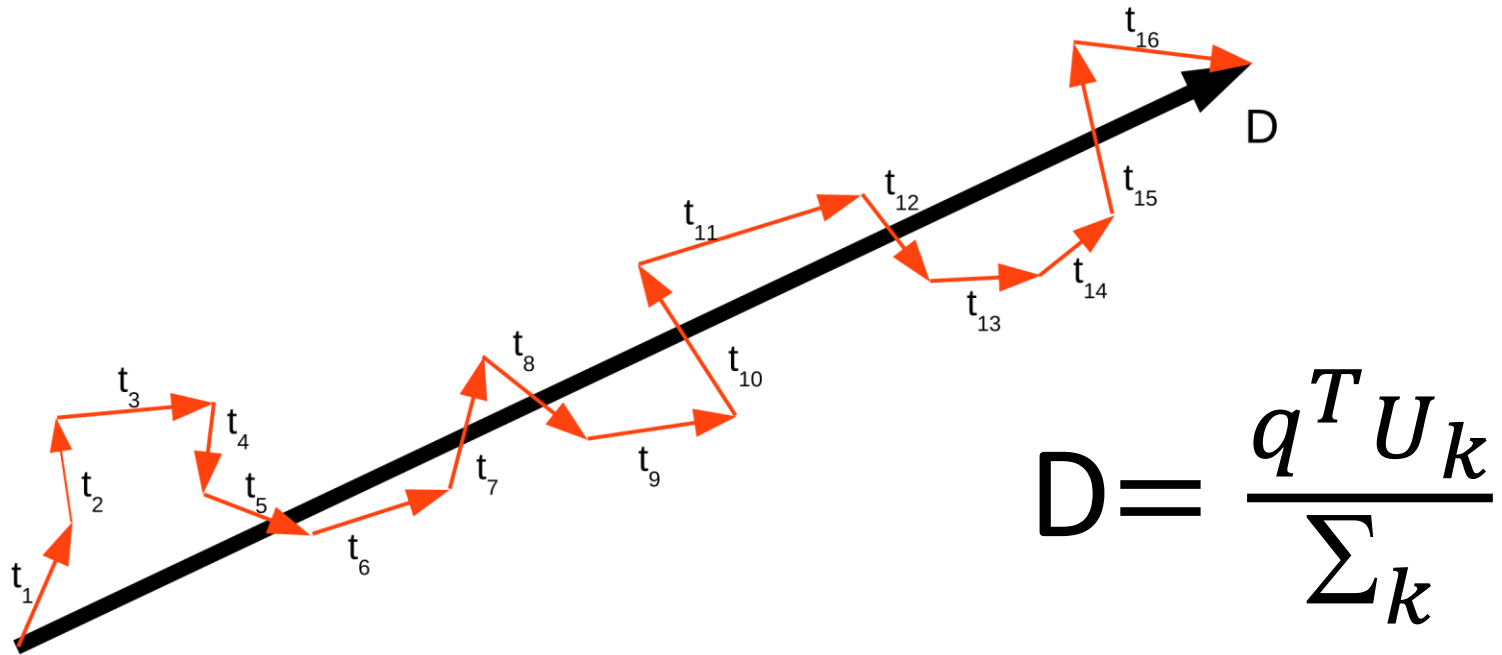
Meaning





Compositionality Constraint

The meaning of a document is the sum of the meaning of its words



Compositionality Constraint Corollary

The meaning of a word is defined by
the documents in which it appears
(and does not appear)

Meaning

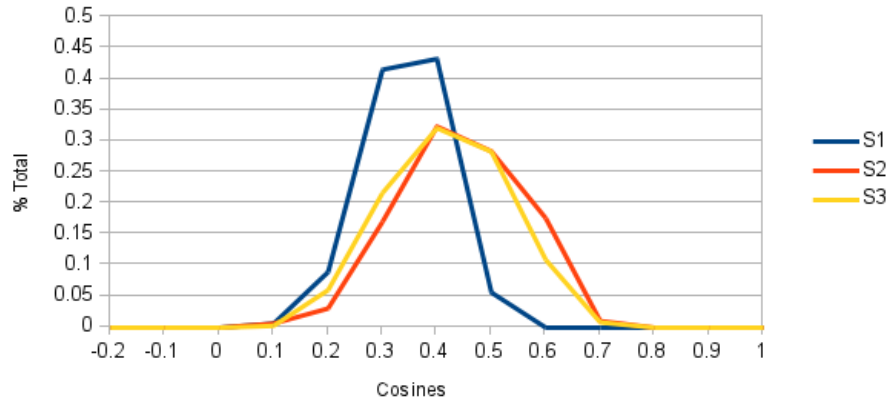
The Mapping system consists of:

- Term Vector Dictionary
- Singular Values

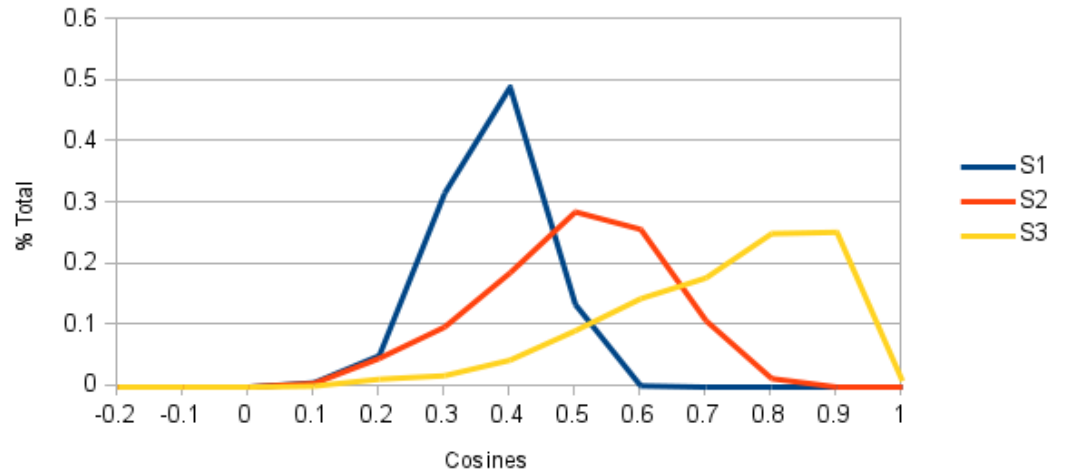


Motivation

Documents to Document Centroid



Driving Documents



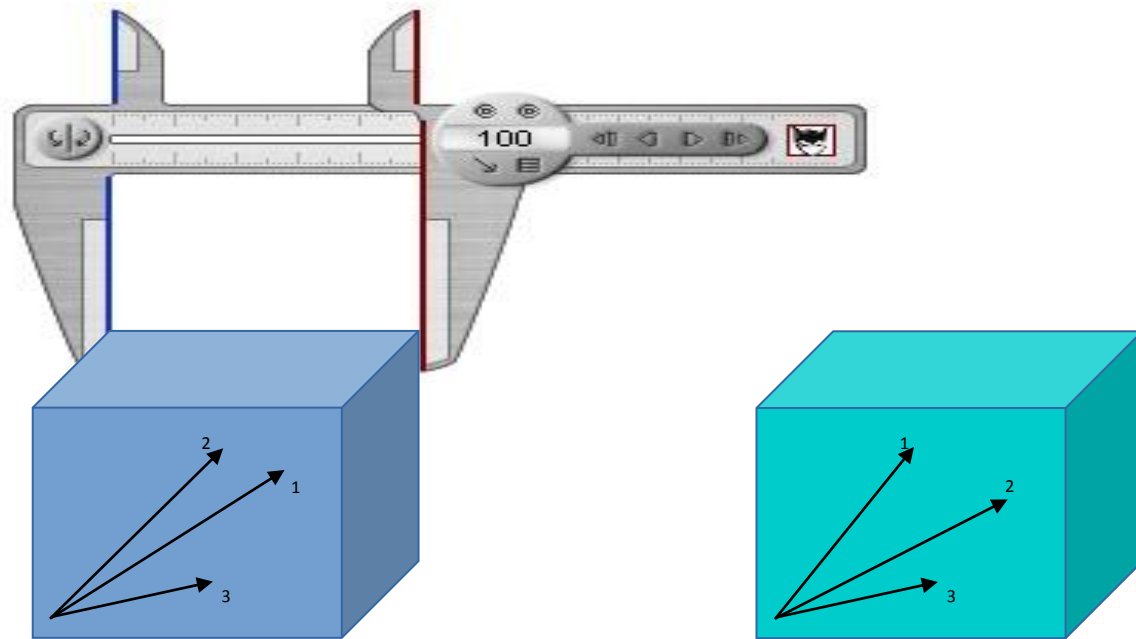
Objective

Find a measure or set of measures that can quantify the difference between two spaces

Measures

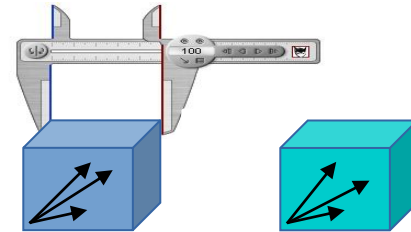
- Direct Comparison
- Projected Content Comparison
- Rotated Item Comparison

Direct Comparison Measures

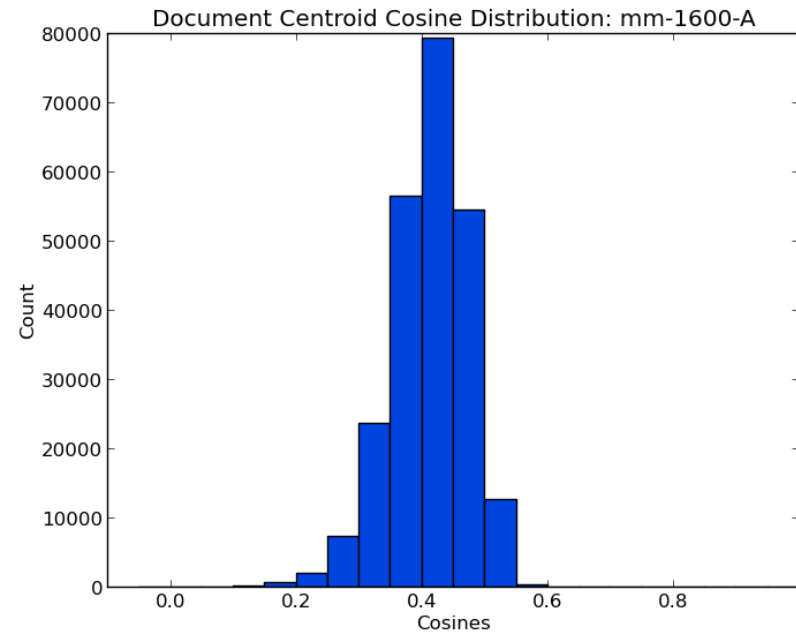
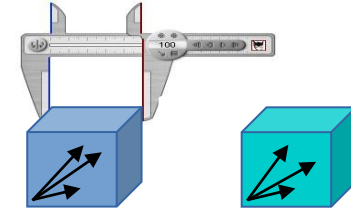
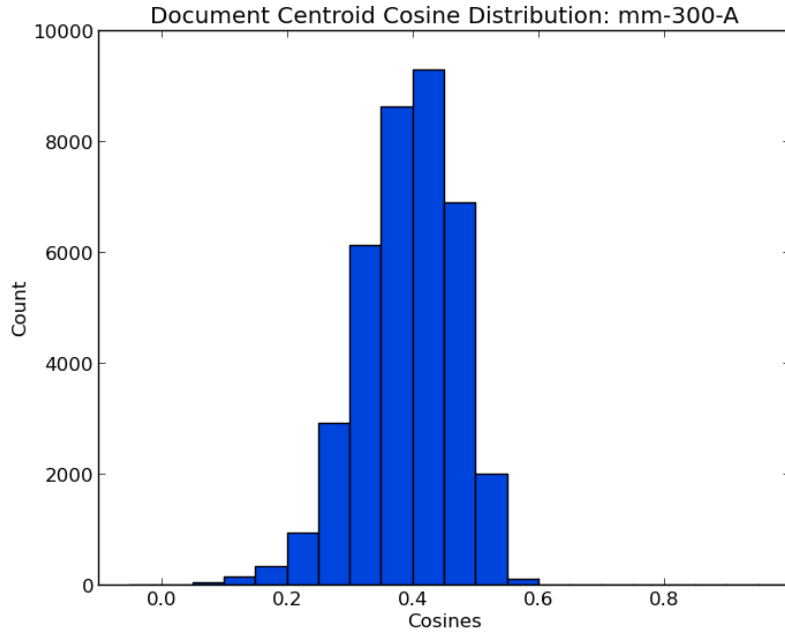


Individual Space Measures

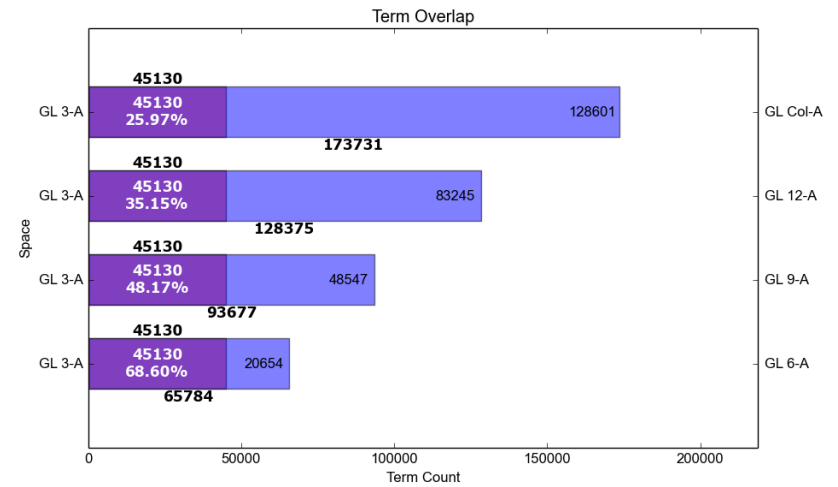
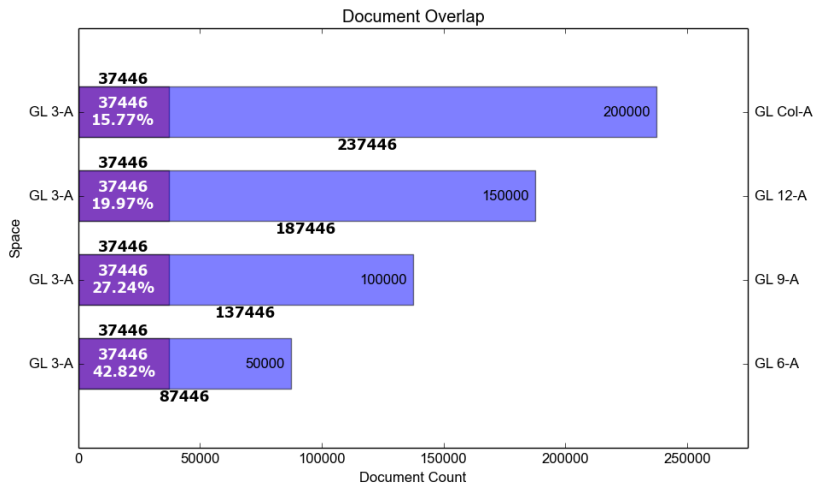
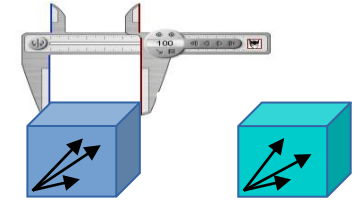
- Document Count
- Term Count
- Non-zeroes



Distribution Analysis

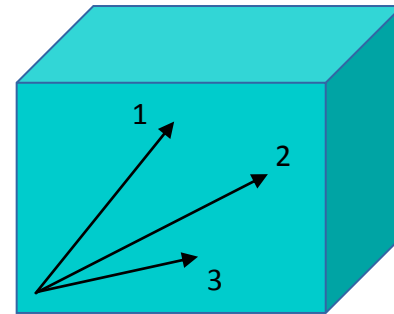
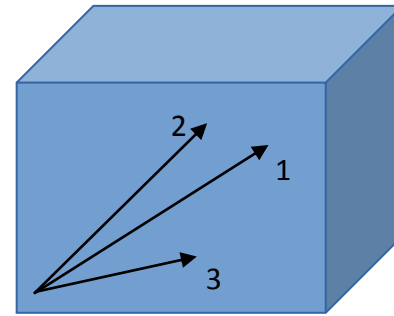


Term and Document Overlap

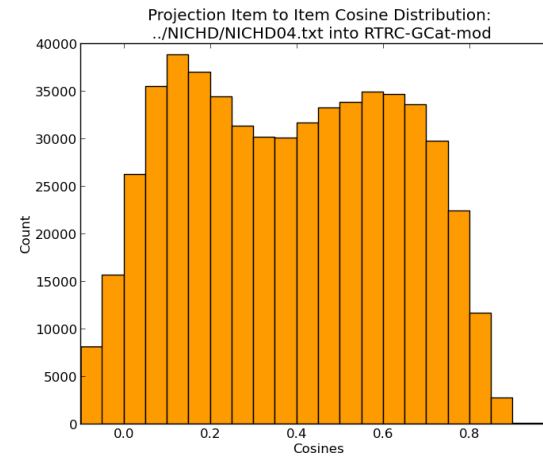
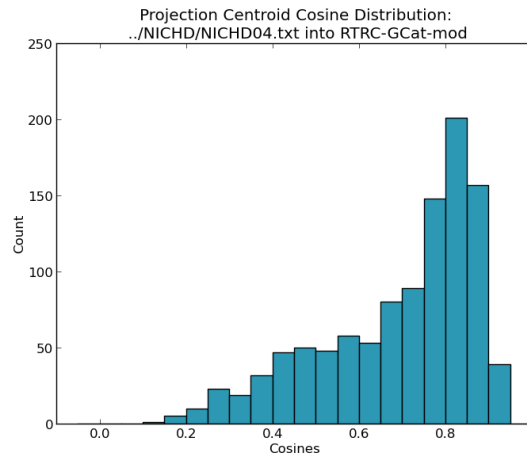
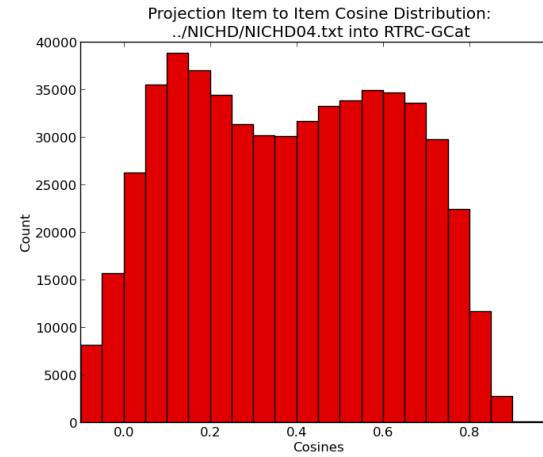
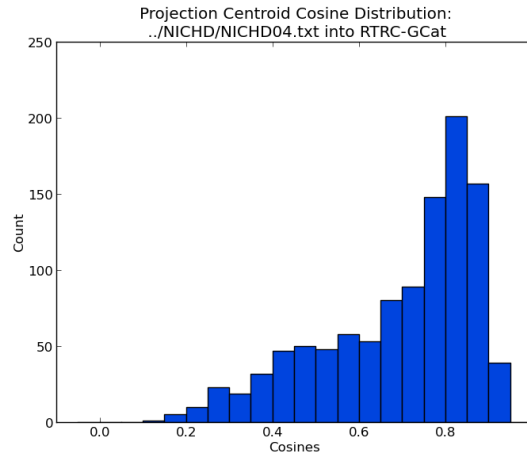
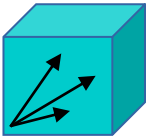
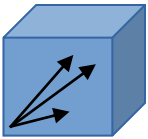


Projected Content Comparisons

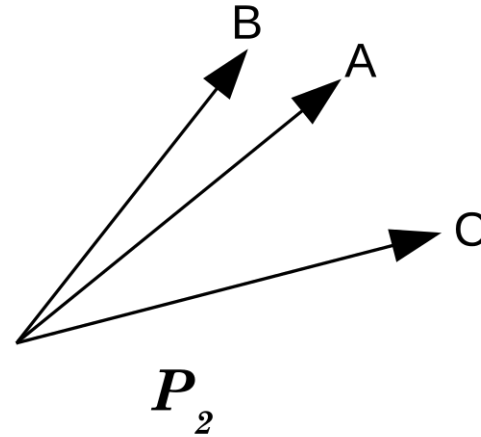
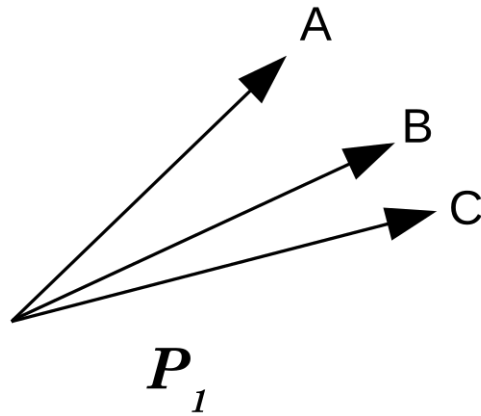
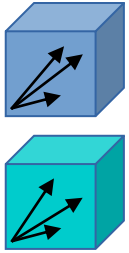
Matched items
projected into each
space



Projected Item Distribution

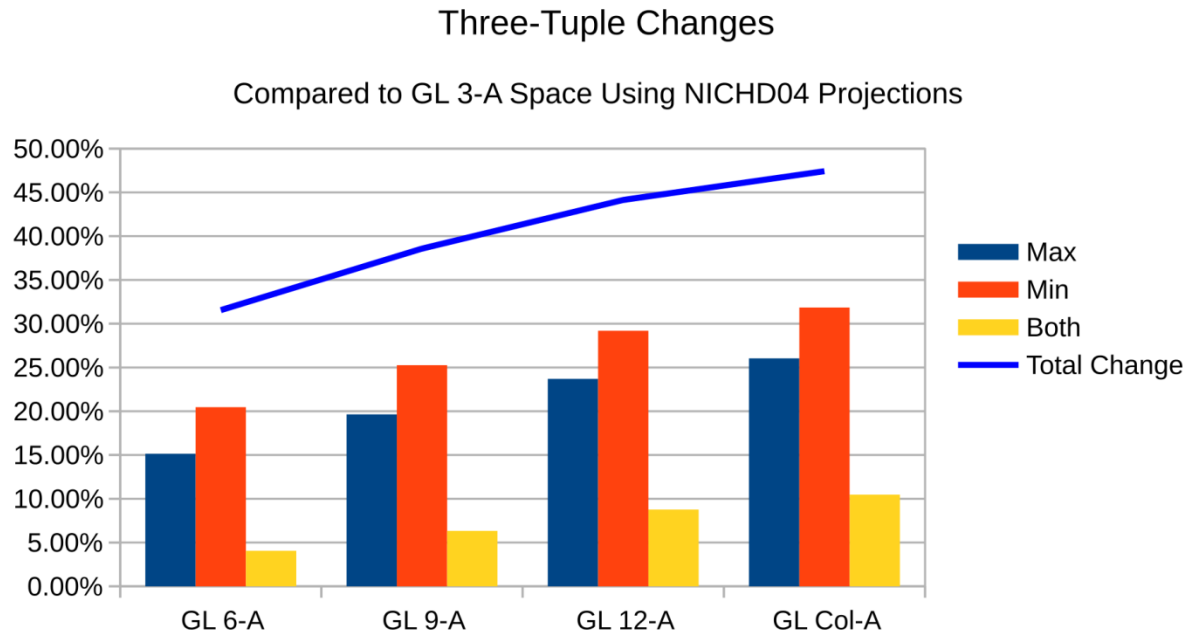
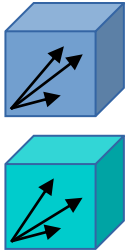


Three-Tuple Comparisons

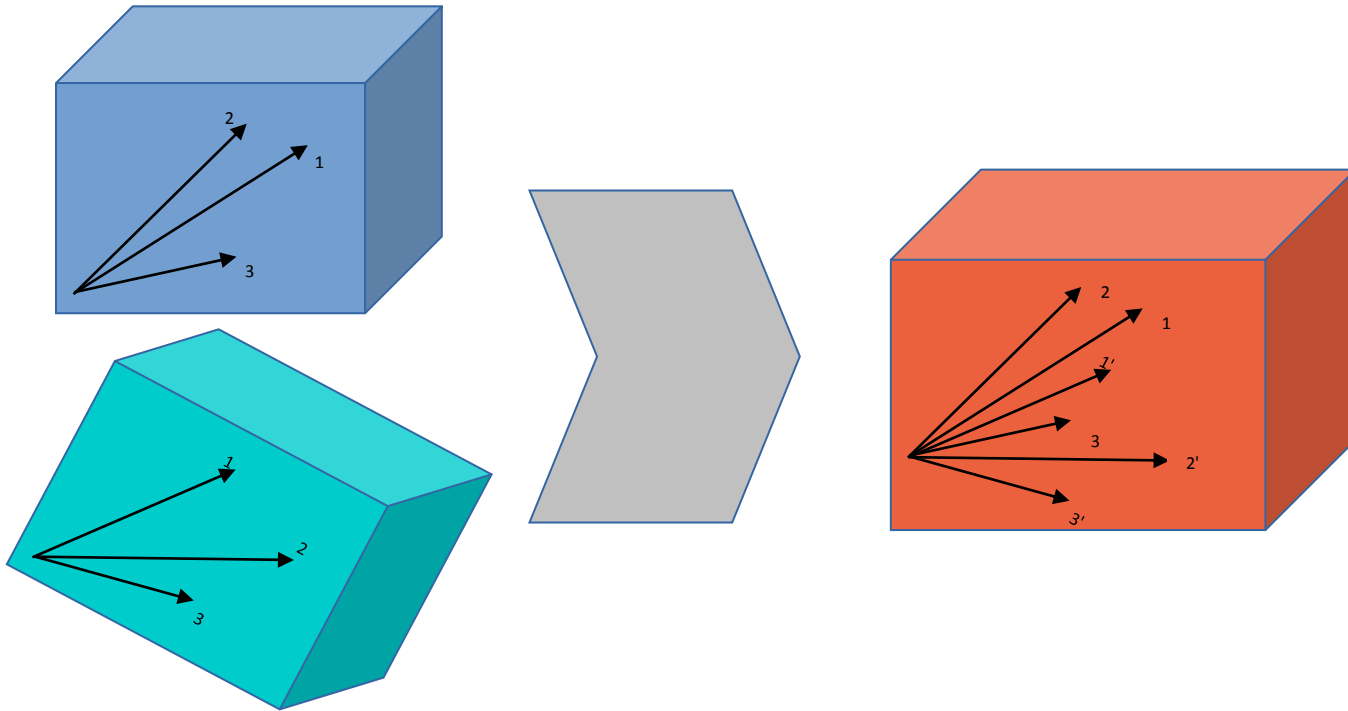


$$\{(A, B, C) \mid A = p_i, B = p_j, C = p_k, \text{ where } i \neq j \neq k, \forall p \in P\}$$

Three-Tuple Relationship Changes



Rotations and Transform Comparisons



The Transform

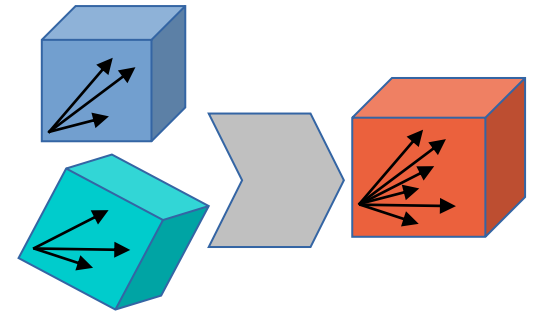
$$A_1 = \text{Project}(A, S_1)$$

$$A_2 = \text{Project}(A, S_2)$$

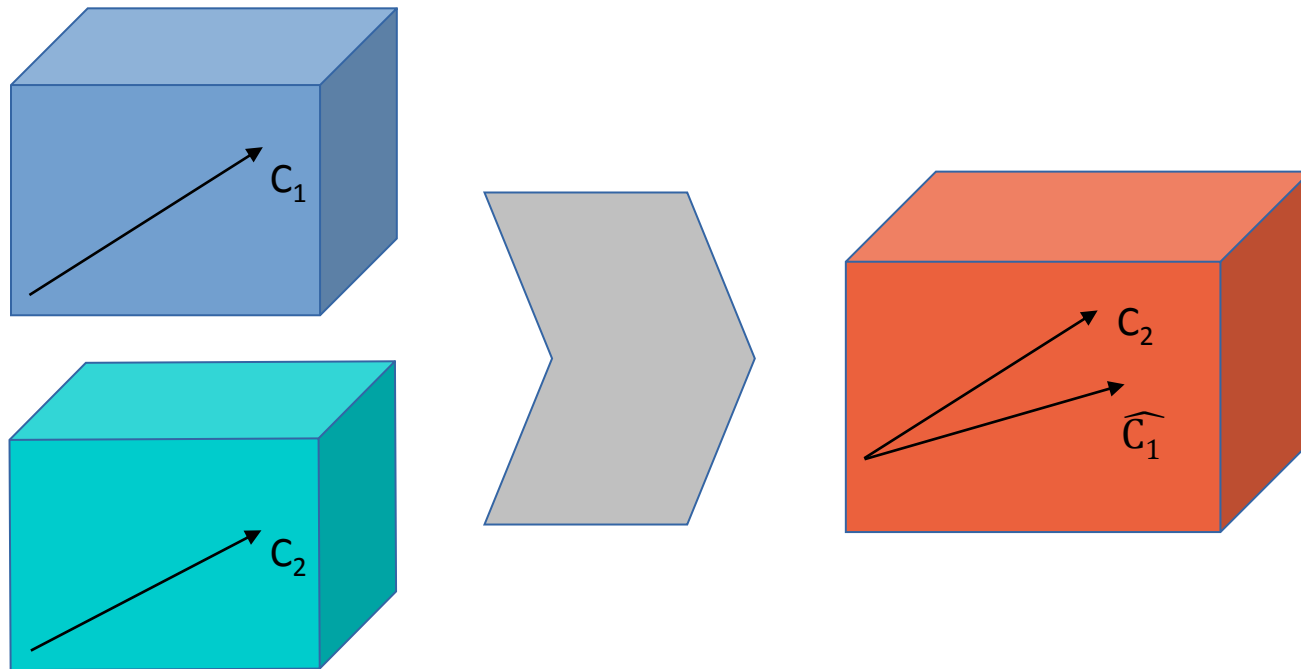
$$A_1^T A_2 = U \Sigma V^T$$

$$Q = UV^T$$

$$\|A_1 Q - A_2\|_F$$



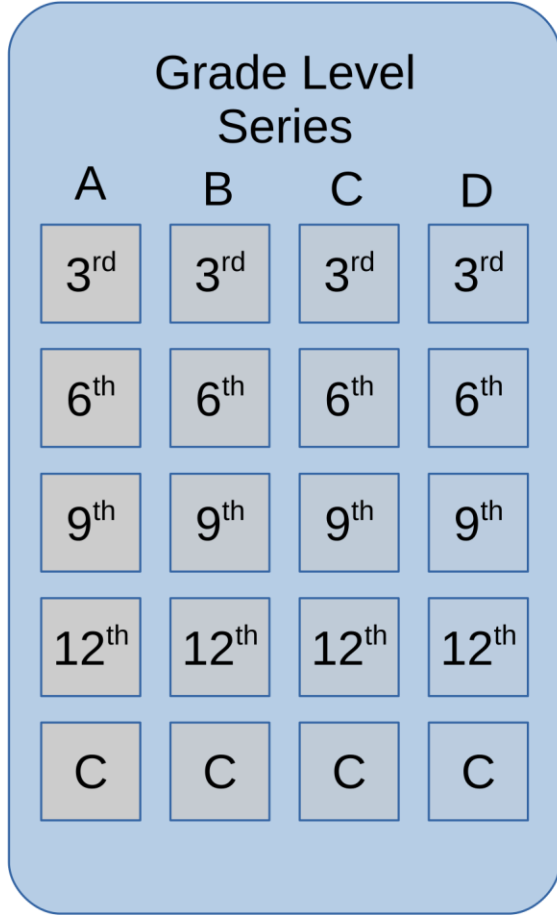
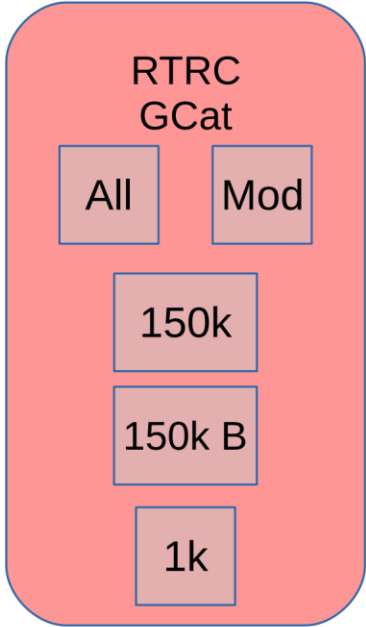
Comparative Space Centroid Analysis



Overlapping Term Vector Norm

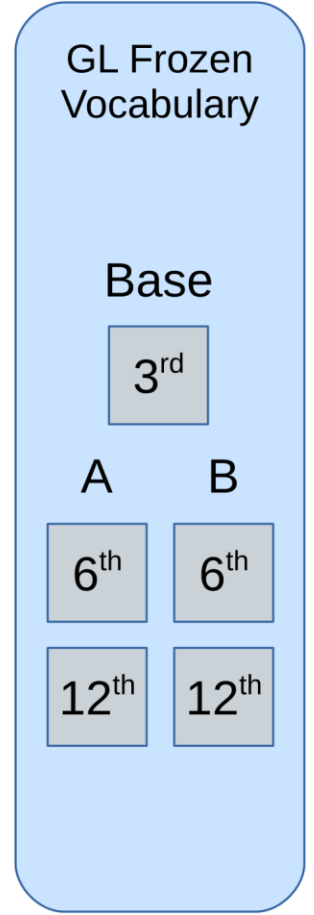
$$\|T_1 Q - T_2\|_F = \|\hat{T}\|_F = \sqrt{\sum_{i=1}^{|\hat{T}|} \sum_{j=1}^k |\hat{t}_{i,j}|^2}$$

6



Large
6th

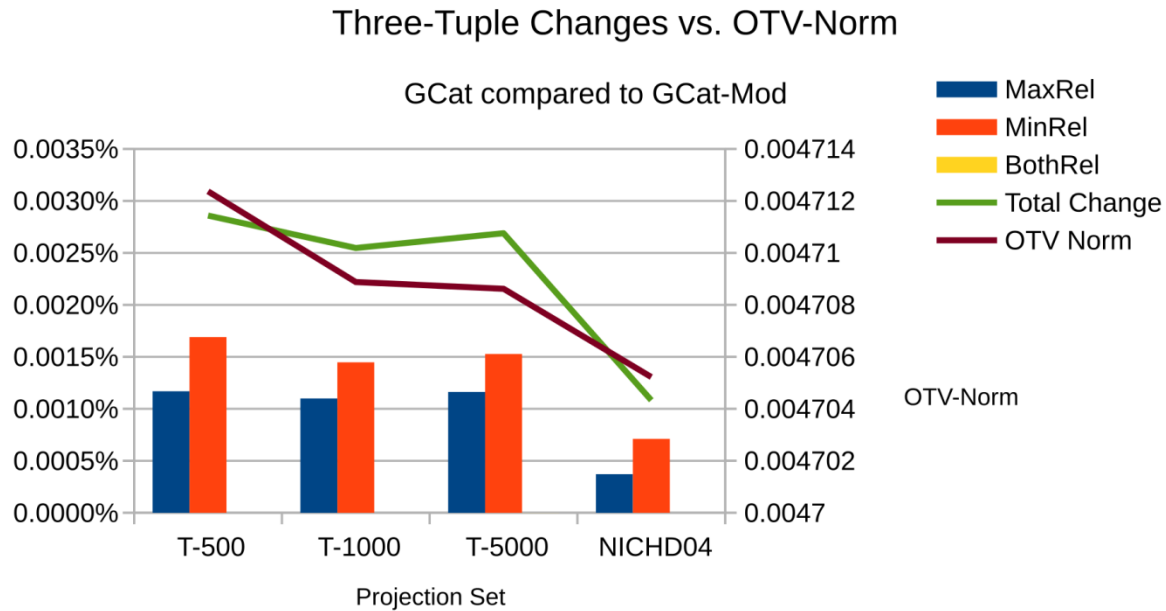
Large
9th



Projection/Anchor Sets

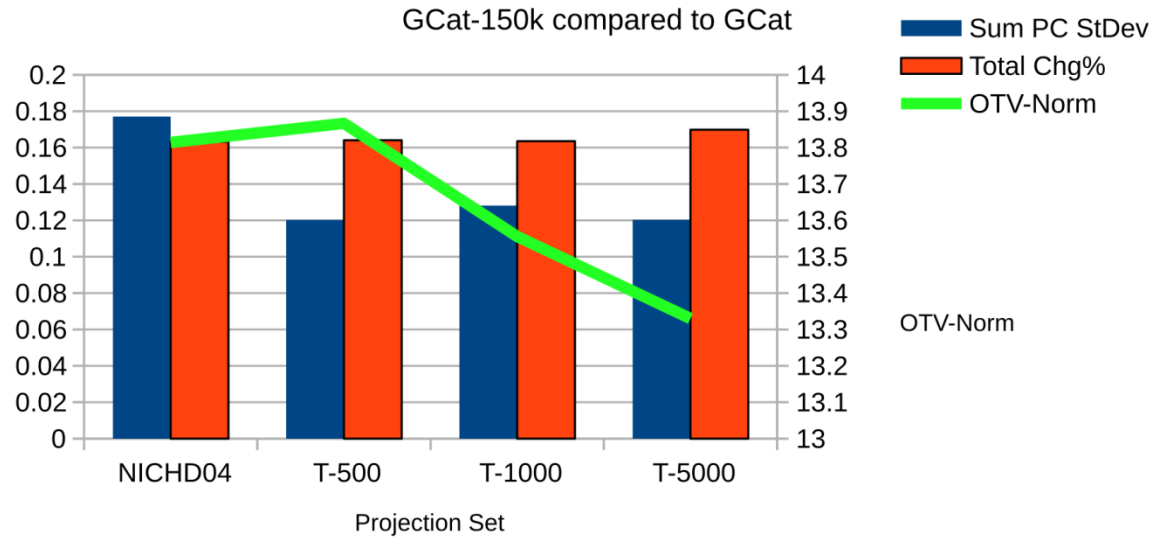
Set	Documents	Unique Terms	Term Instances
NICHD04	1,060	5,912	70,063
T-500	500	16,317	123,668
T-1000	1,000	24,319	252,372
T-5000	5,000	49,995	1,281,749

Control Experiment



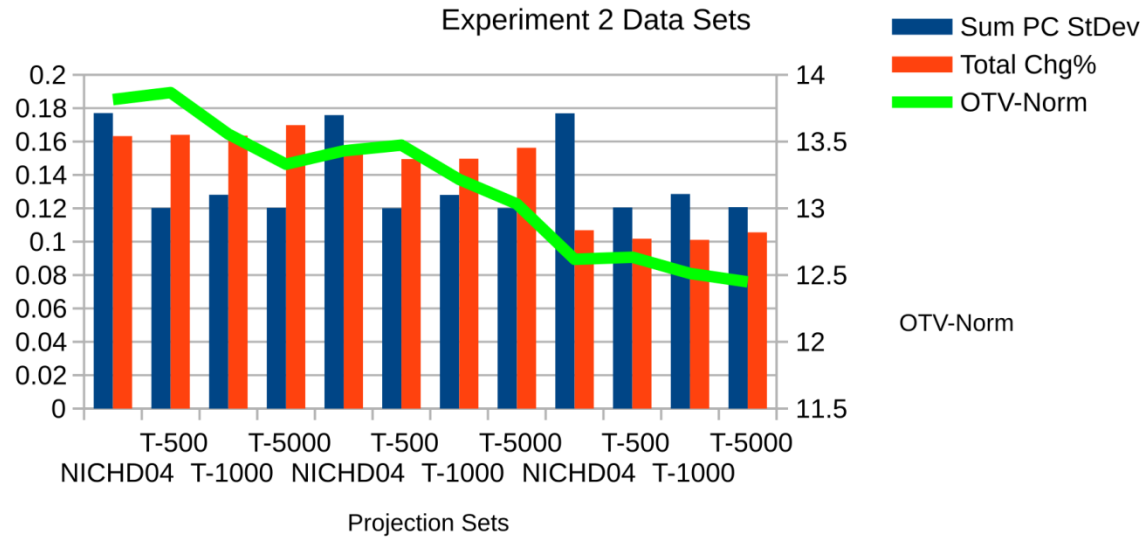
General Experiment

OTV-Norm vs. Total Change% and PC Standard Deviation

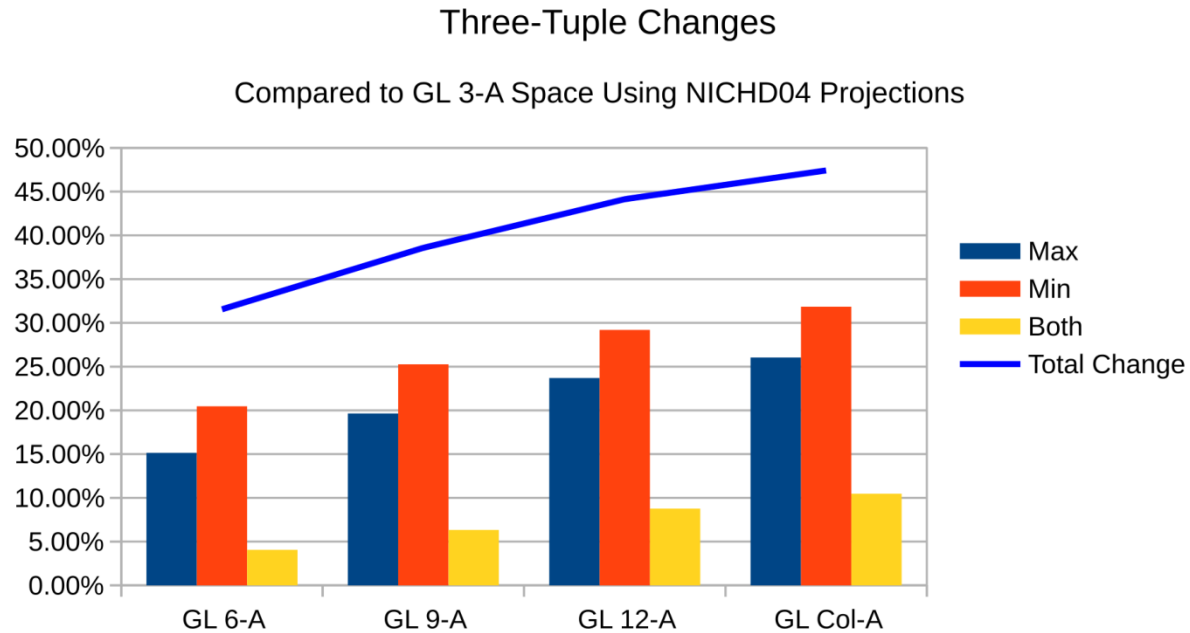


General Experiment

OTV-Norm Vs. Total Change% and Standard Deviation

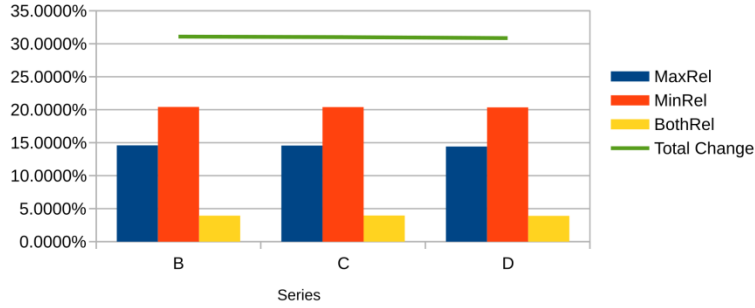


Grade Level Series Experiment



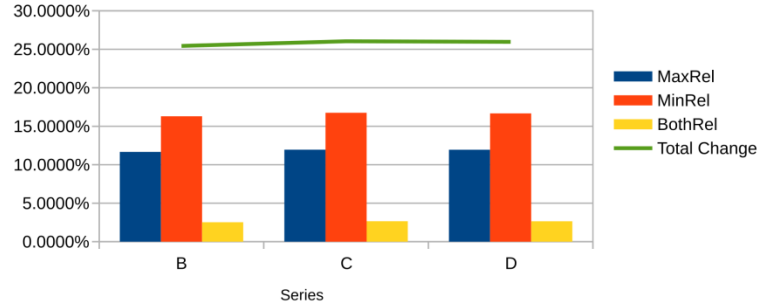
Three-Tuple changes - GL 3-A Cross Series

Using NICH04 Projections



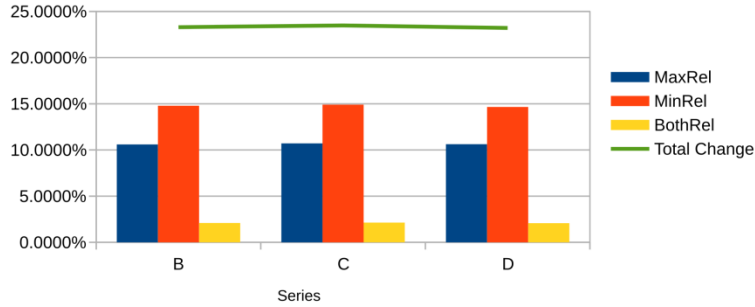
Three-Tuple Changes -GL 6-A Cross Series

Using NICH04 Projections



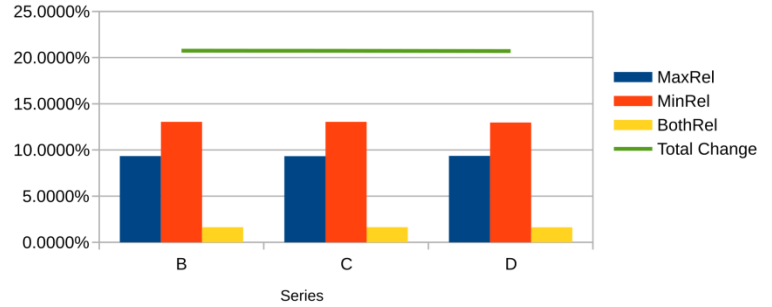
Three-Tuple Changes - GL 9-A Cross Series

Using NICH04 Projections



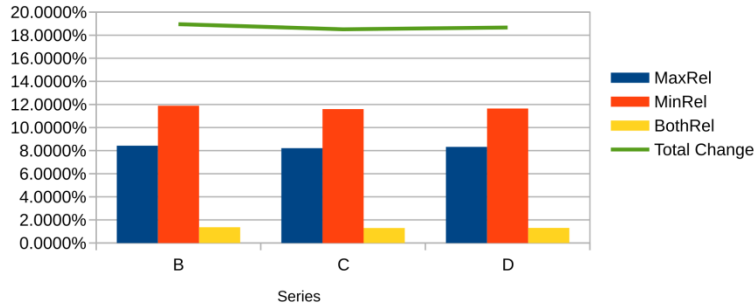
Three-Tuple Changes - GL 12-A Cross Series

Using NICH04 Projections

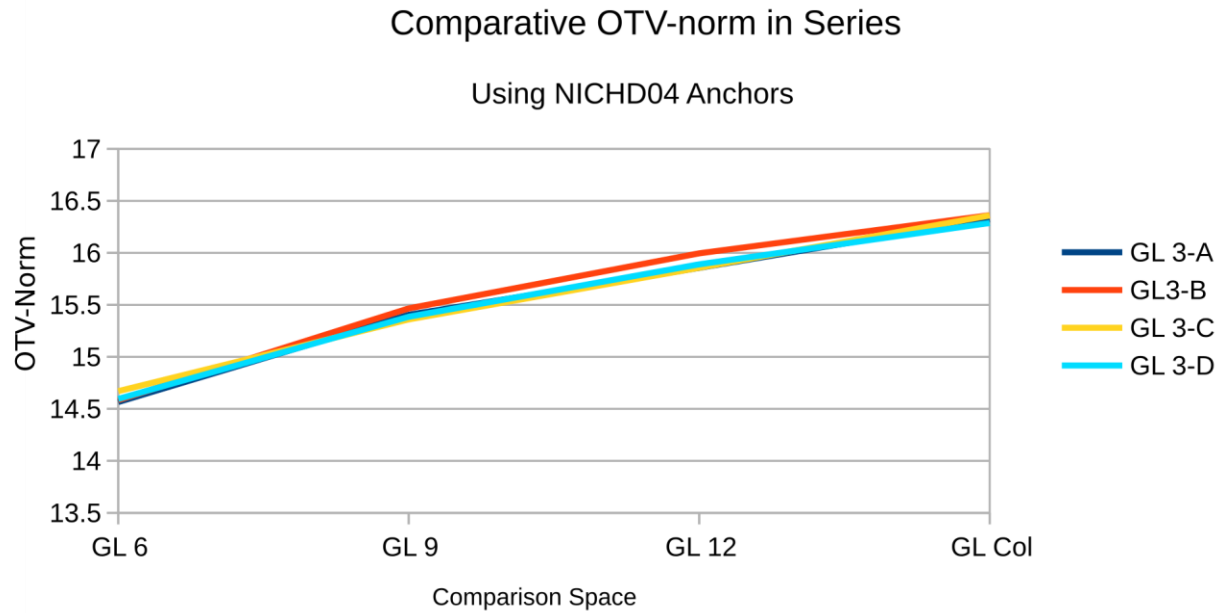


Three-Tuple Changes - GL Col-A Cross Series

Using NICH04 Projections



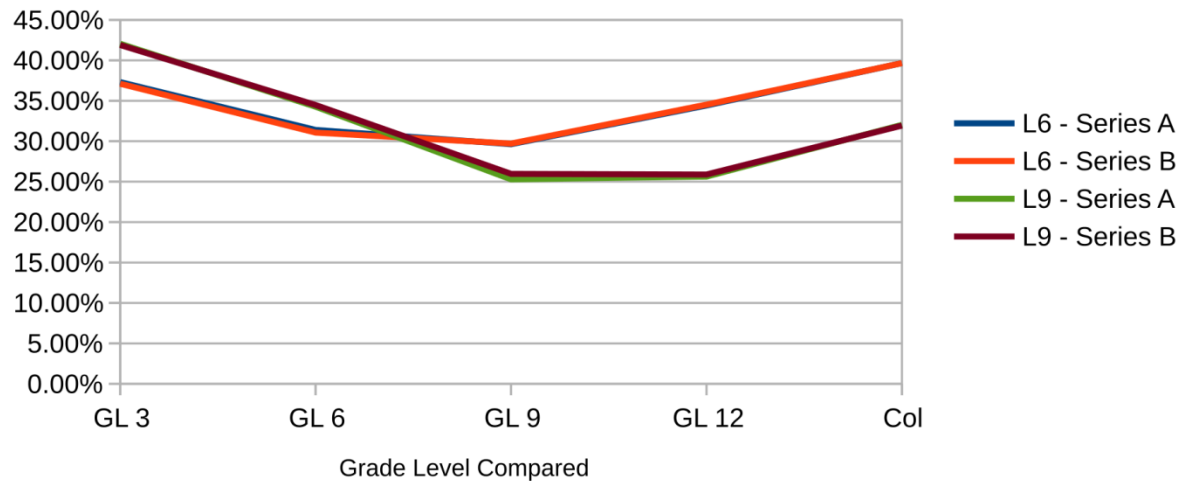
Grade Level Series OTV-Norm



Large Volume Experiment

Large Space Comparisons

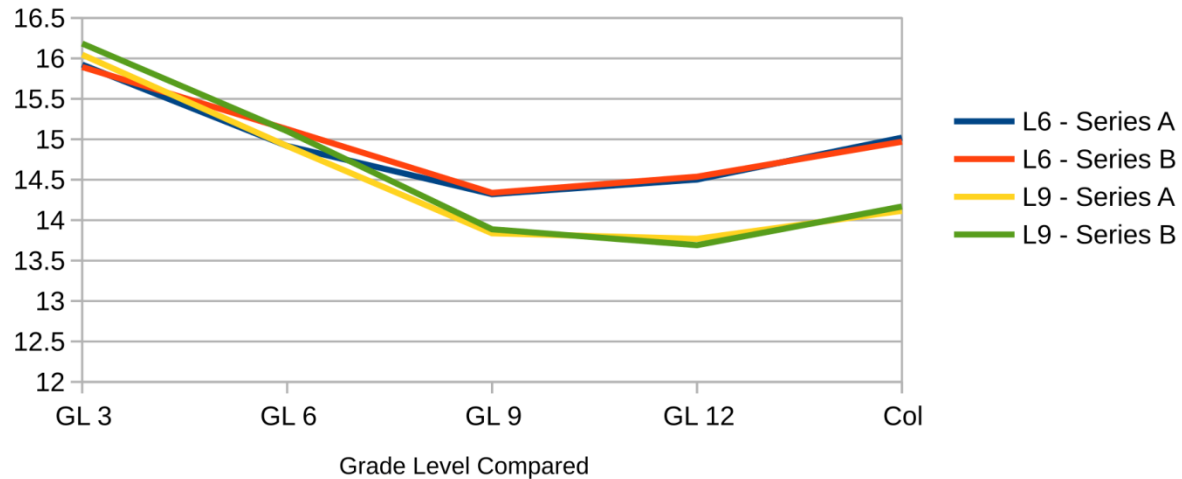
TC% Between Spaces



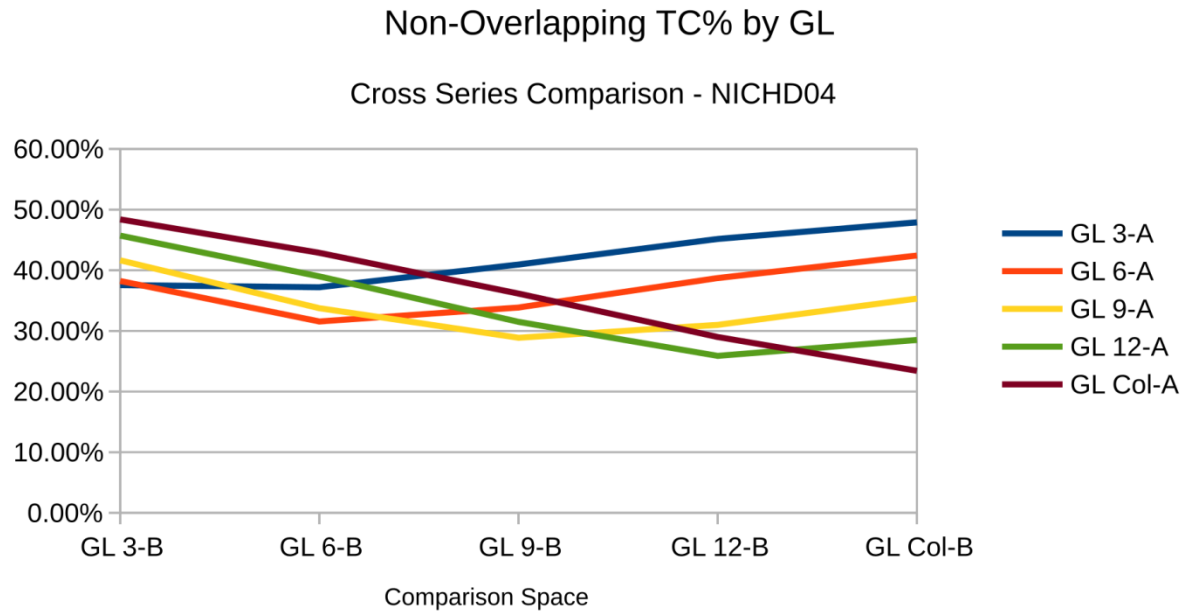
Large Volume Experiment

Large Space Comparisons

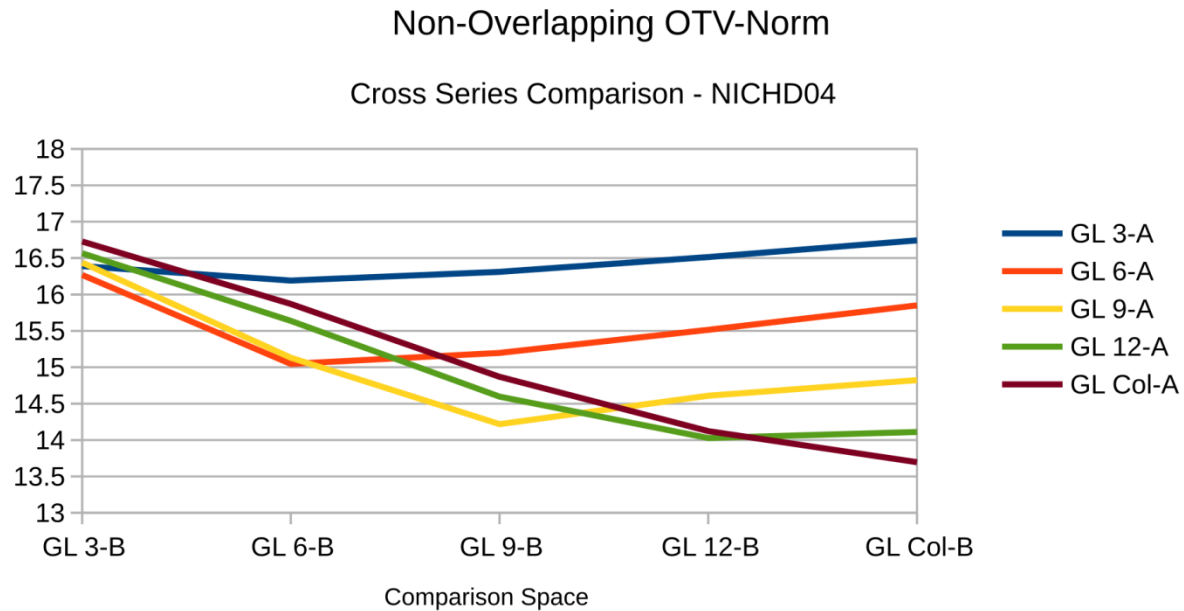
OTV-Norm Between Spaces



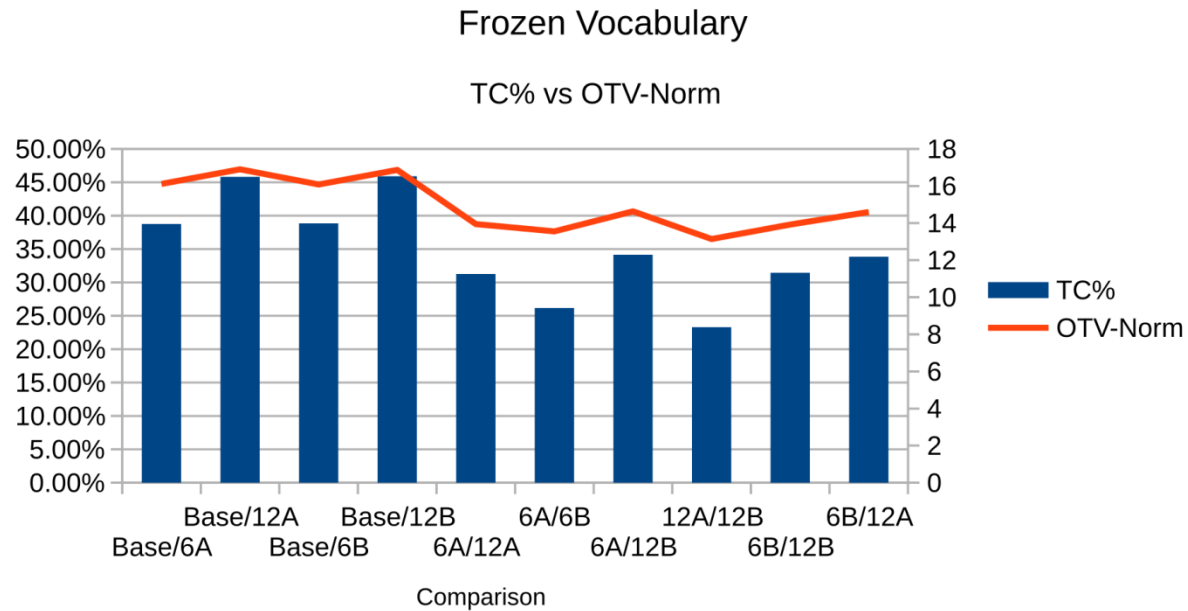
Non-overlapping Series Experiment



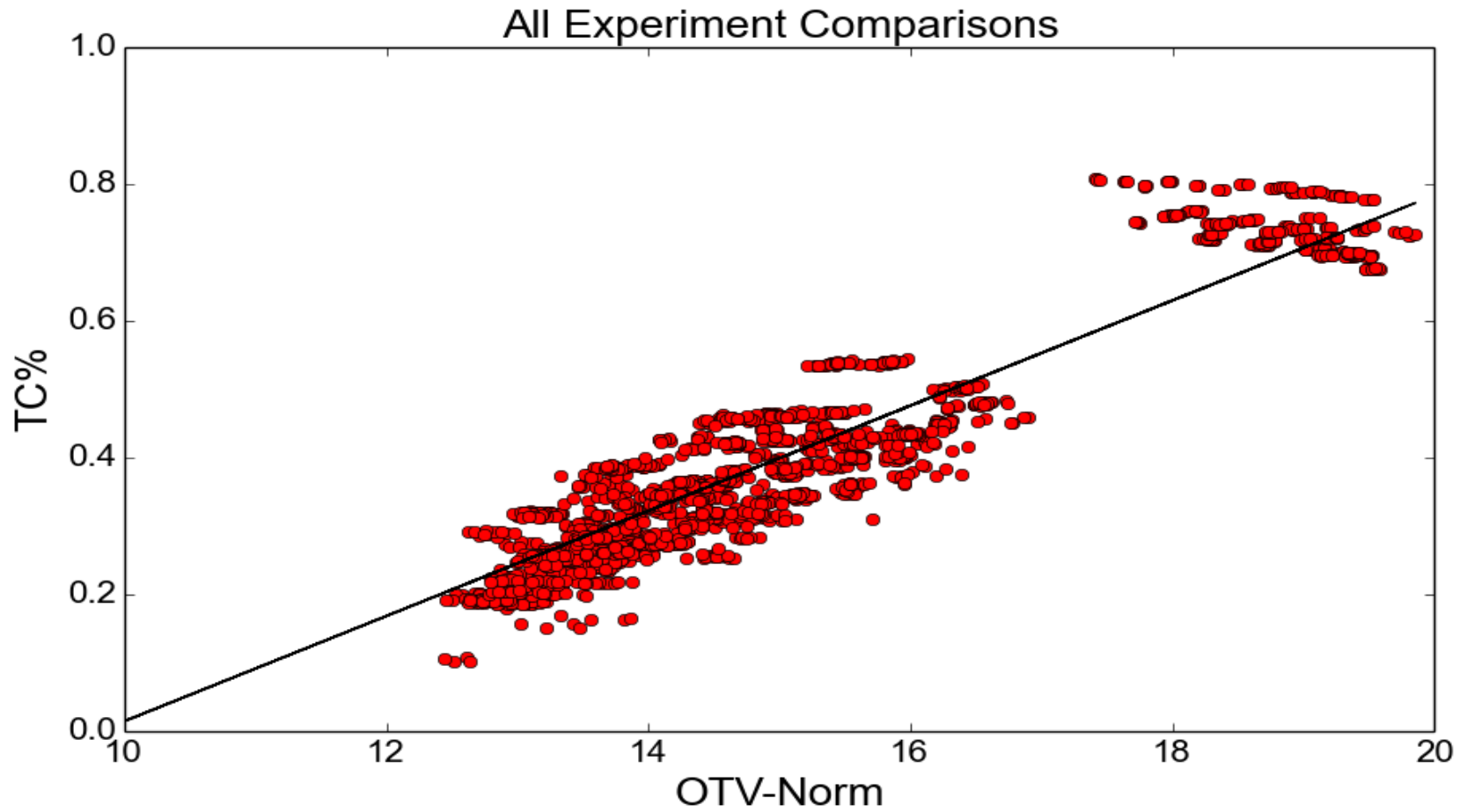
Non-Overlapping Series OTV-Norm



Frozen Vocabulary Experiment



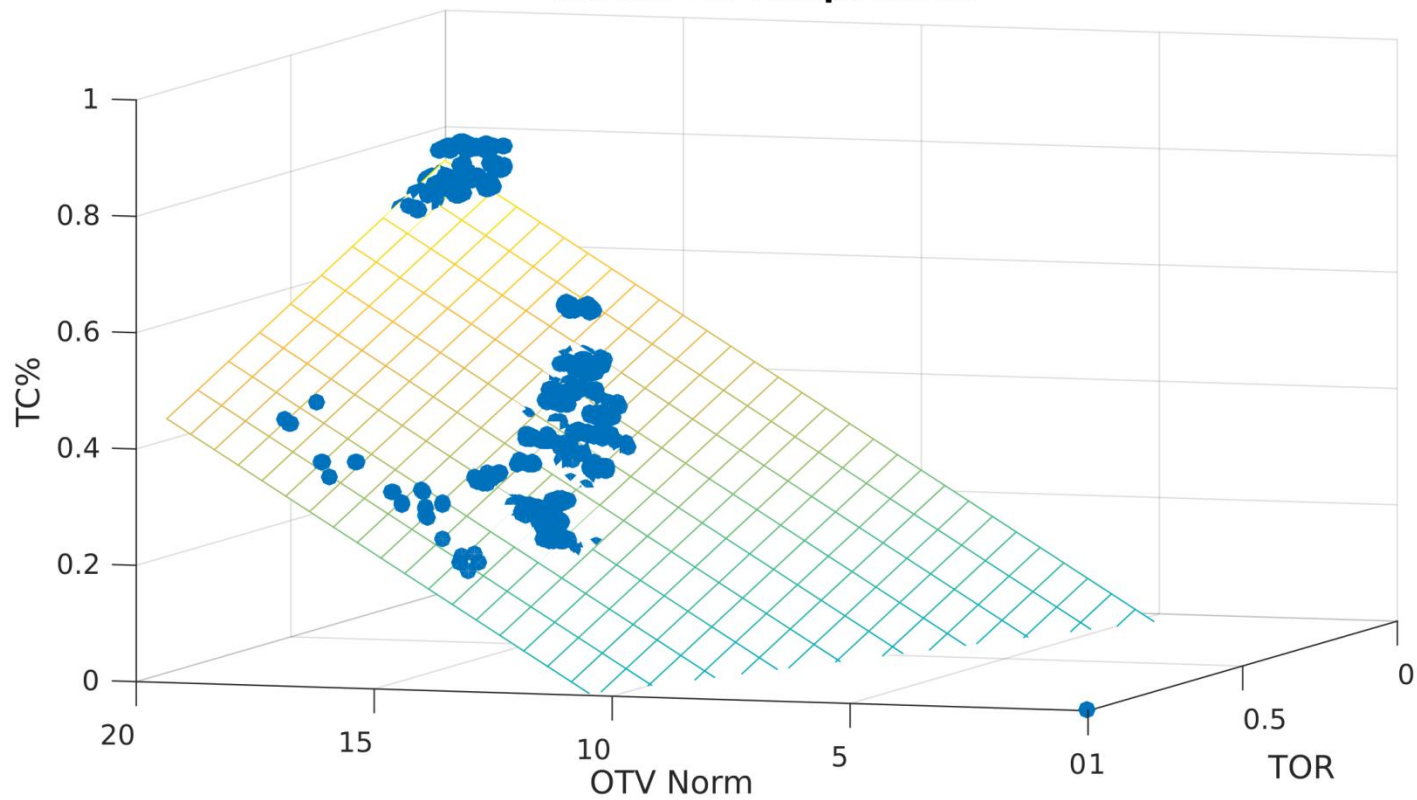
OTV-Norm



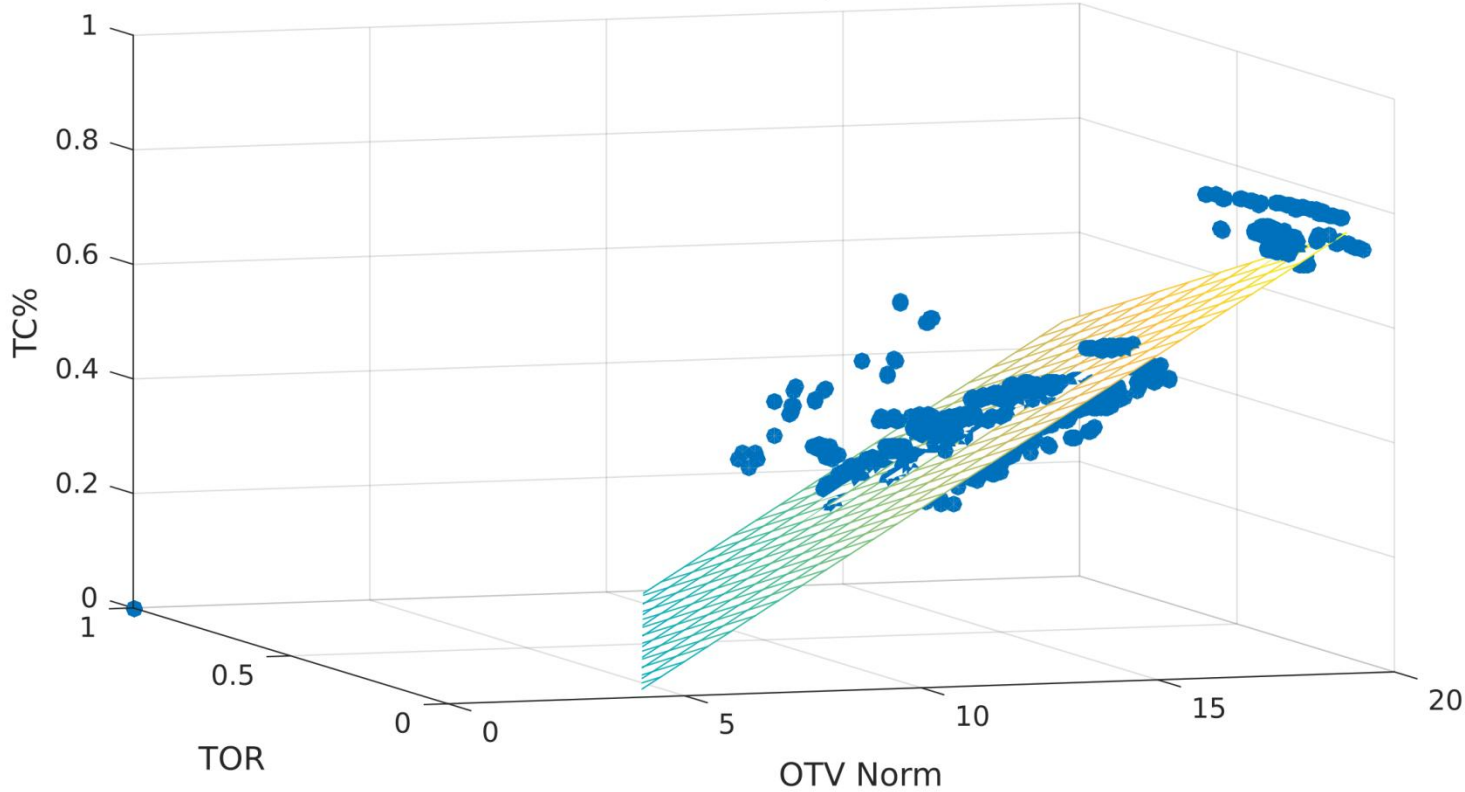
Semantic Measurement Model

$$\begin{aligned} TC\% &\approx -0.207882 \\ &+ 0.0507194(OTVNorm) \\ &+ -0.339339(TOR) \end{aligned}$$

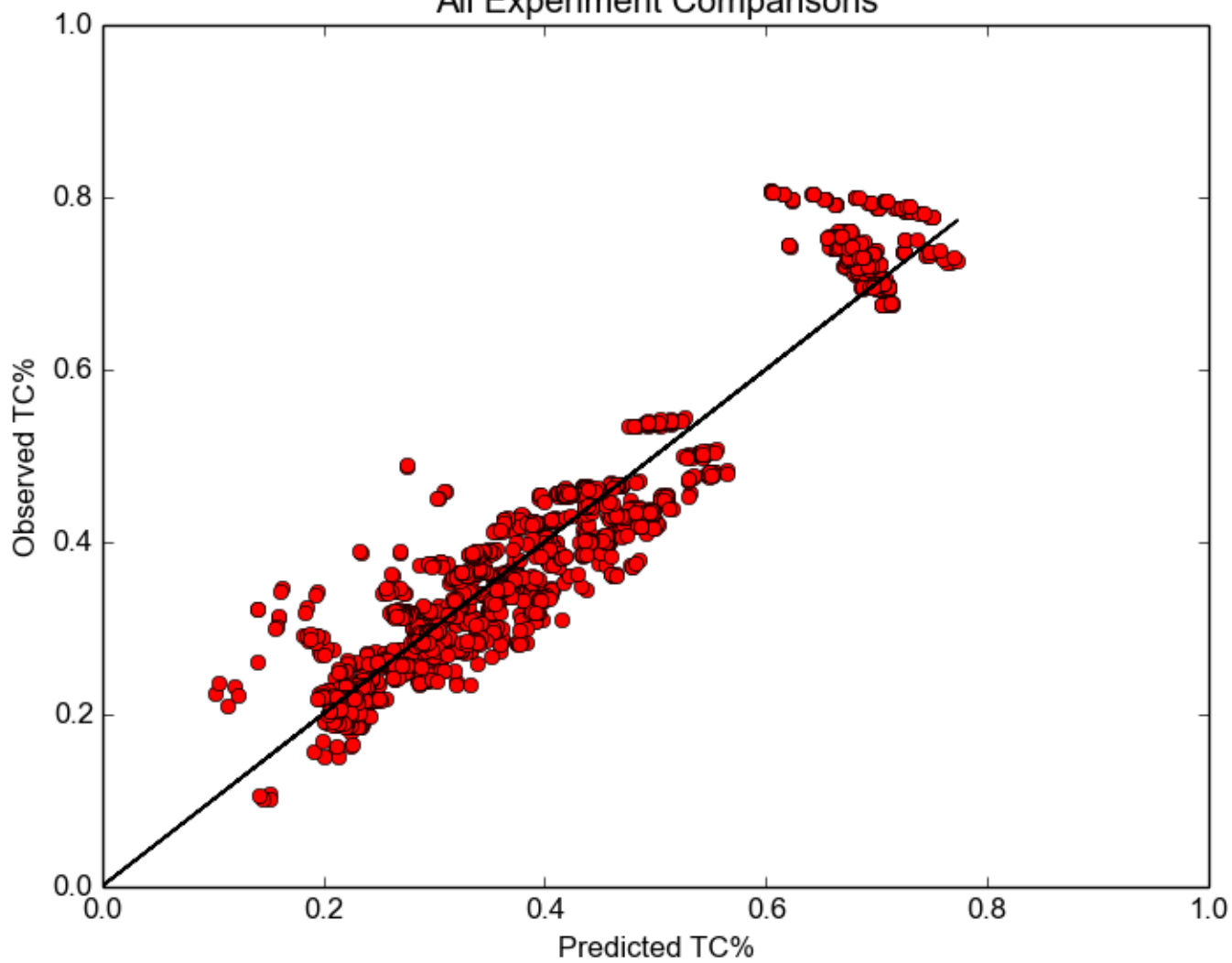
Linear Regression of TOR and OTV-Norm Across All Comparisons



Linear Regression of TOR and OTV-Norm Across All Comparisons



All Experiment Comparisons



Summary of Contributions

- Semantic differences are observable
 - Measurable
 - Quality based
- Similarity not dependent on overlapping content
- OTV-Norm & Semantic Measurement Model
 - Whole-space measurement

Further Research

- Refine the model
 - Anchor set selection/influence
 - Account for non-overlapping terms
 - Investigate non-linear model
- Other questions raised

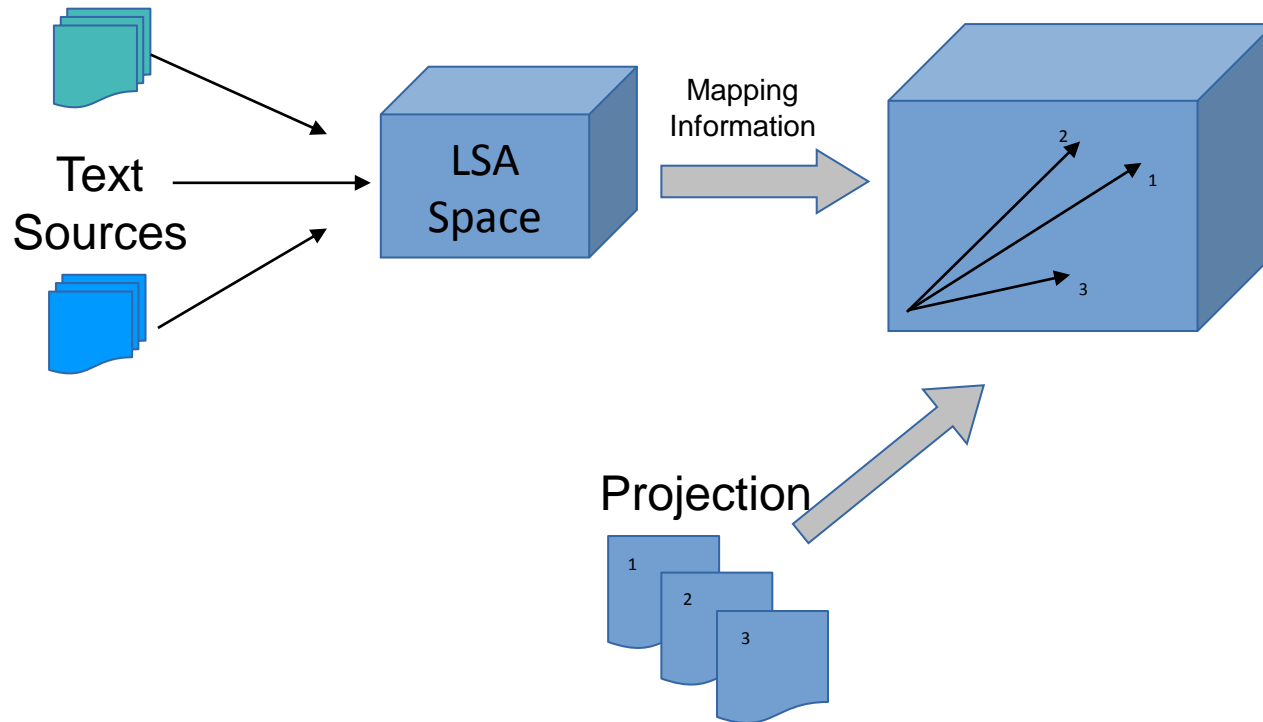
Leverage for Answering Other Questions

- Is it possible to identify key documents that affect the meaning of a space?
- Do additional items added to a space have any impact?
- Is there a point at which adding any items to a space makes no difference?
- Is it possible to identify necessary knowledge that would align two spaces?

Q&A

Backup Slides

Projection of New Content

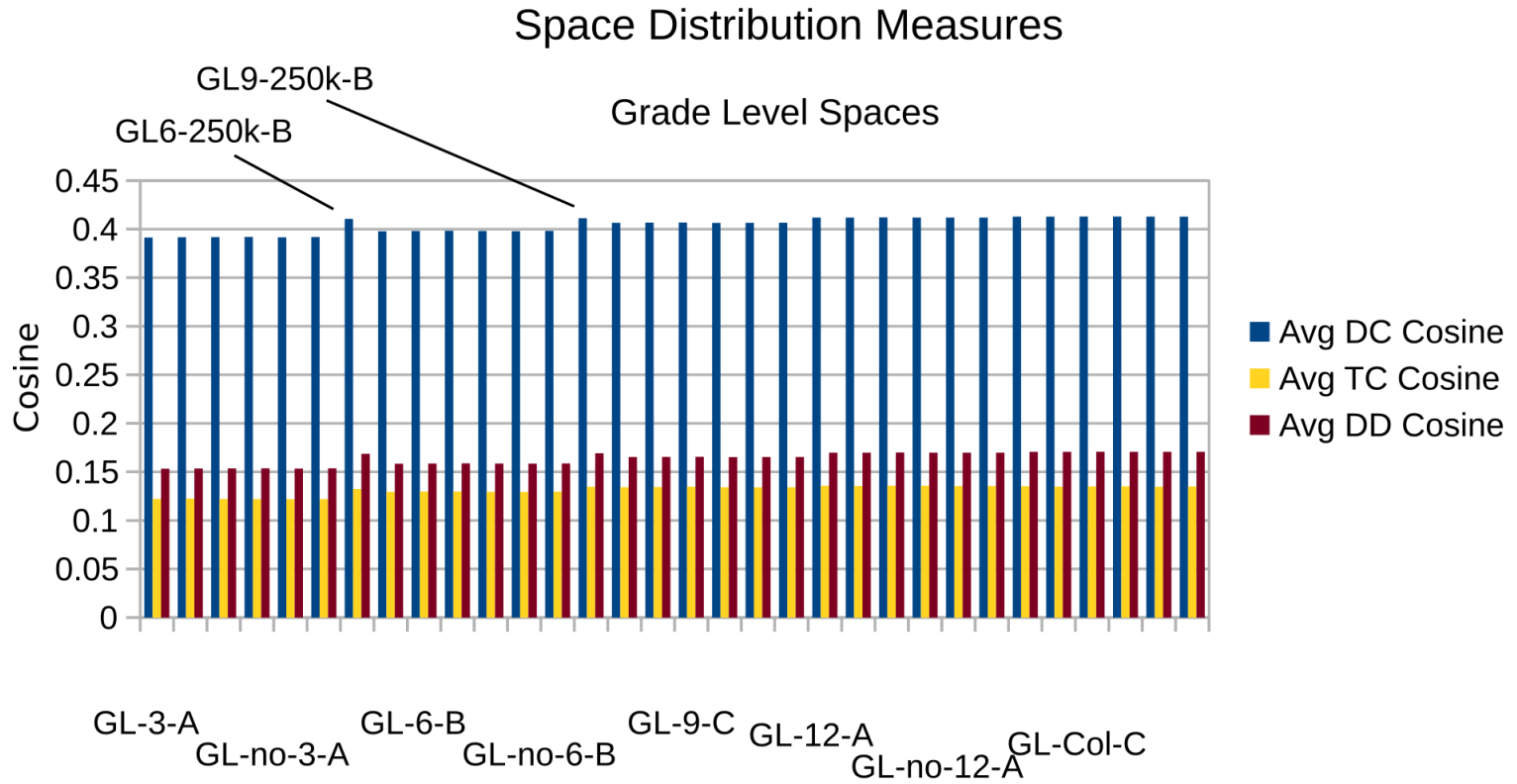


Data

- 42 Spaces
- 592 Comparisons
- 4 Projection Sets
- 4 Anchor Sets
- 26 Measures

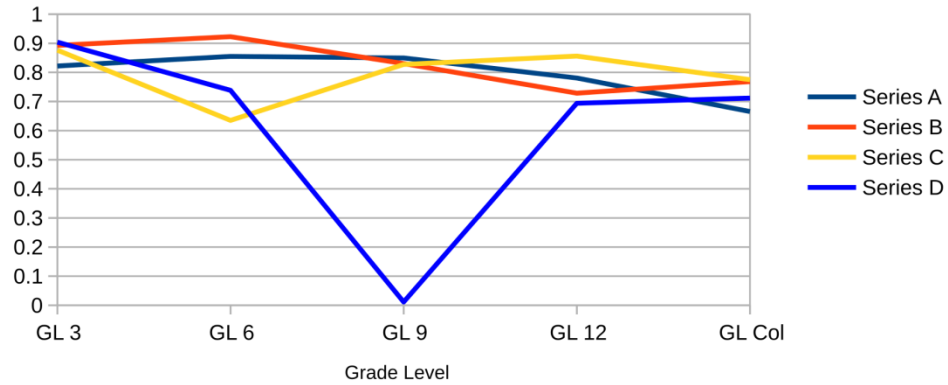
61,568 Data Items Collected

Distribution Analysis



Document Centroid Cosine Comparison

RTRC-GCat Compared to Grade Level Series Using NICHHD04 Anchors



Comparative Term Centroid Cosines in Series

Using NICHHD04 Anchors

