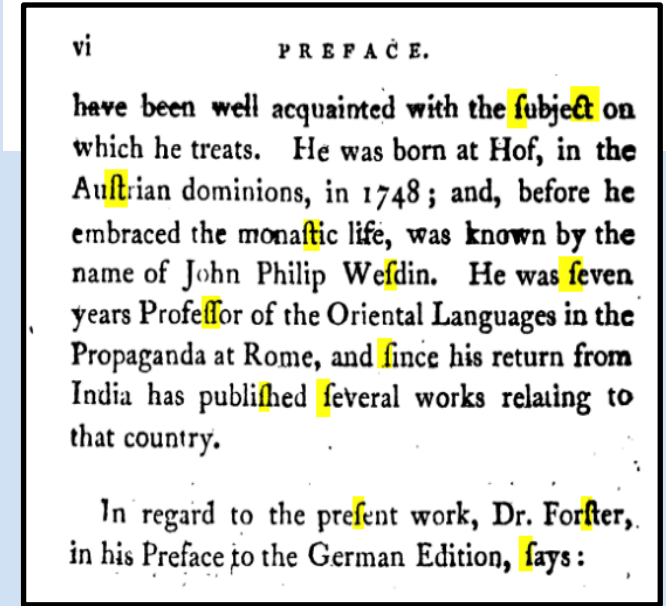
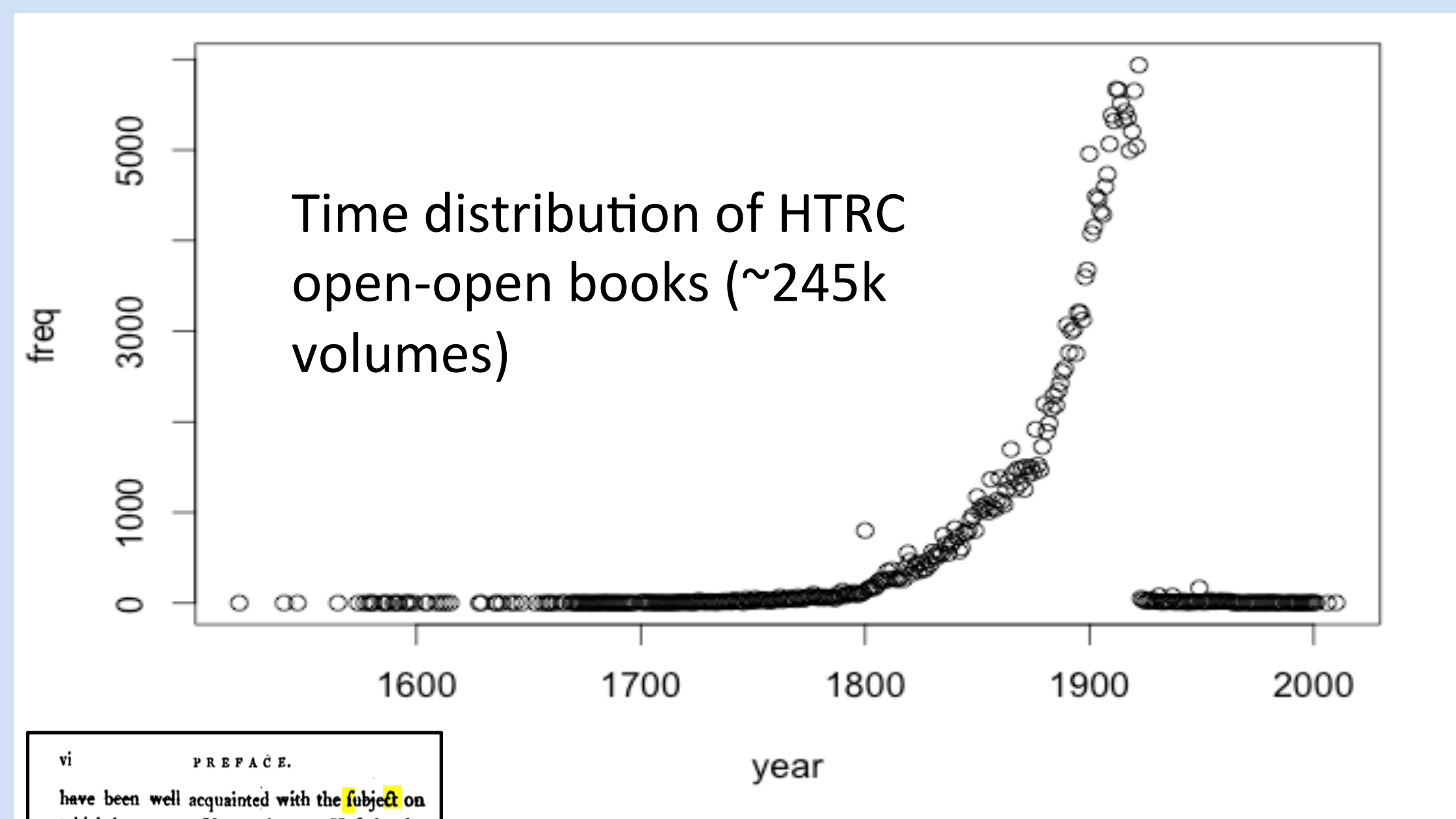


Toward Enhanced Metadata Quality of Large-Scale Digital Libraries: Estimating Volume Time Range

Siyuan Guo, Trevor Edelblute, Bin Dai, Miao Chen, Xiaozhong Liu
Indiana University

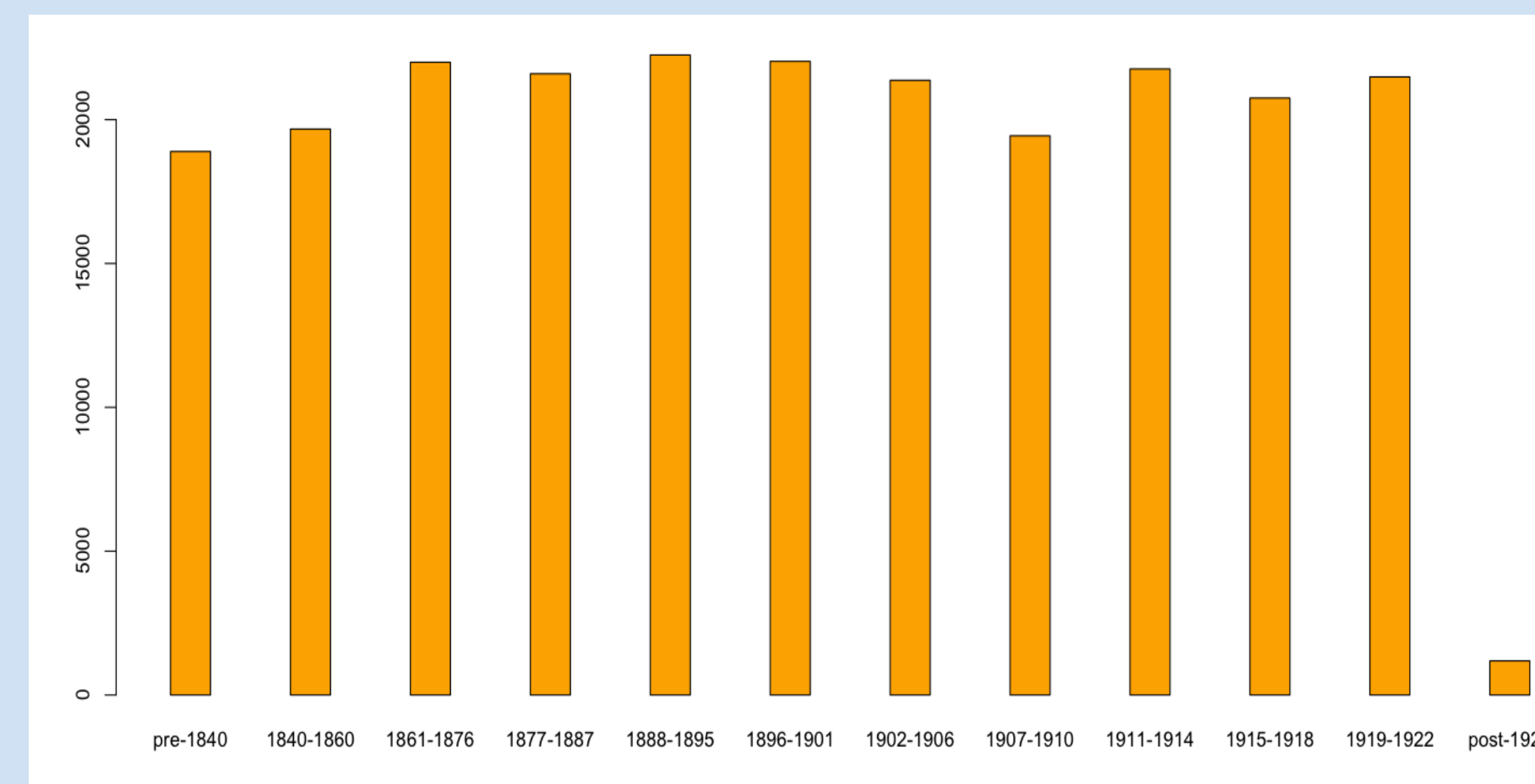
The Big Challenge

It is common that metadata records are incomplete in large-scale digital libraries. E.g. 13% publish dates are missing in HT's non-Google digitized public domain collection



Questions

When is a book published?
Can we tell it by looking at the book content?

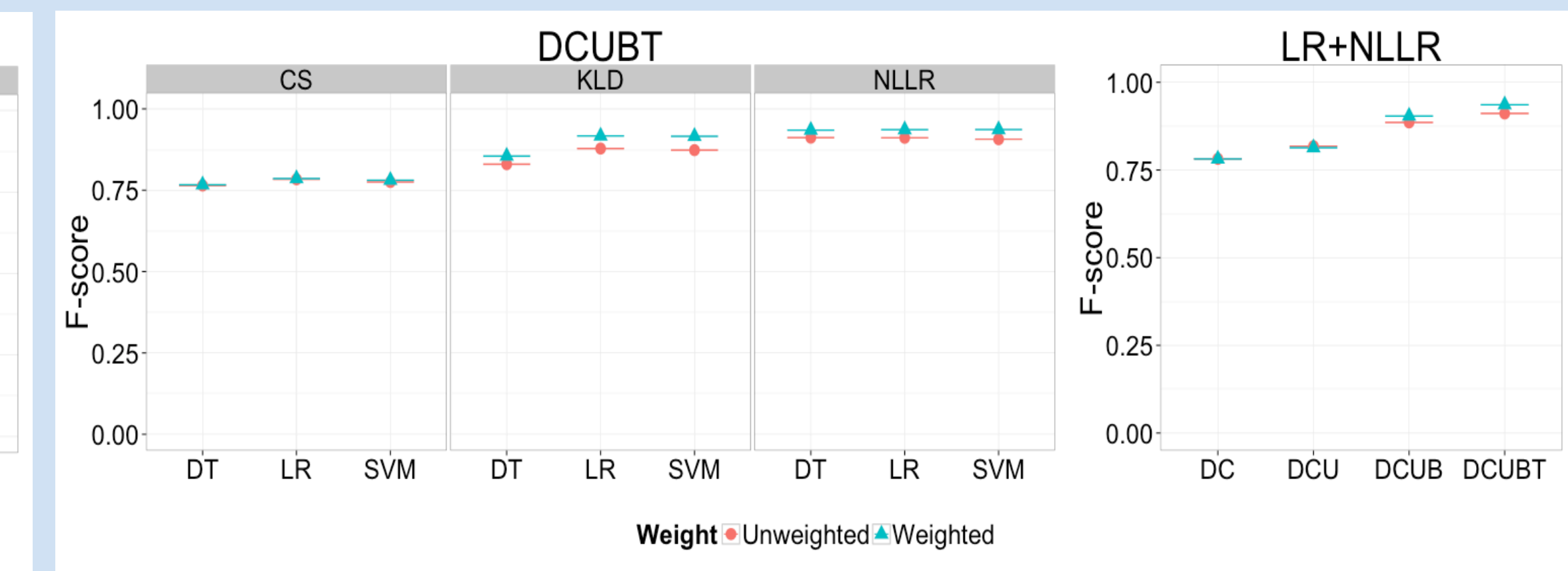
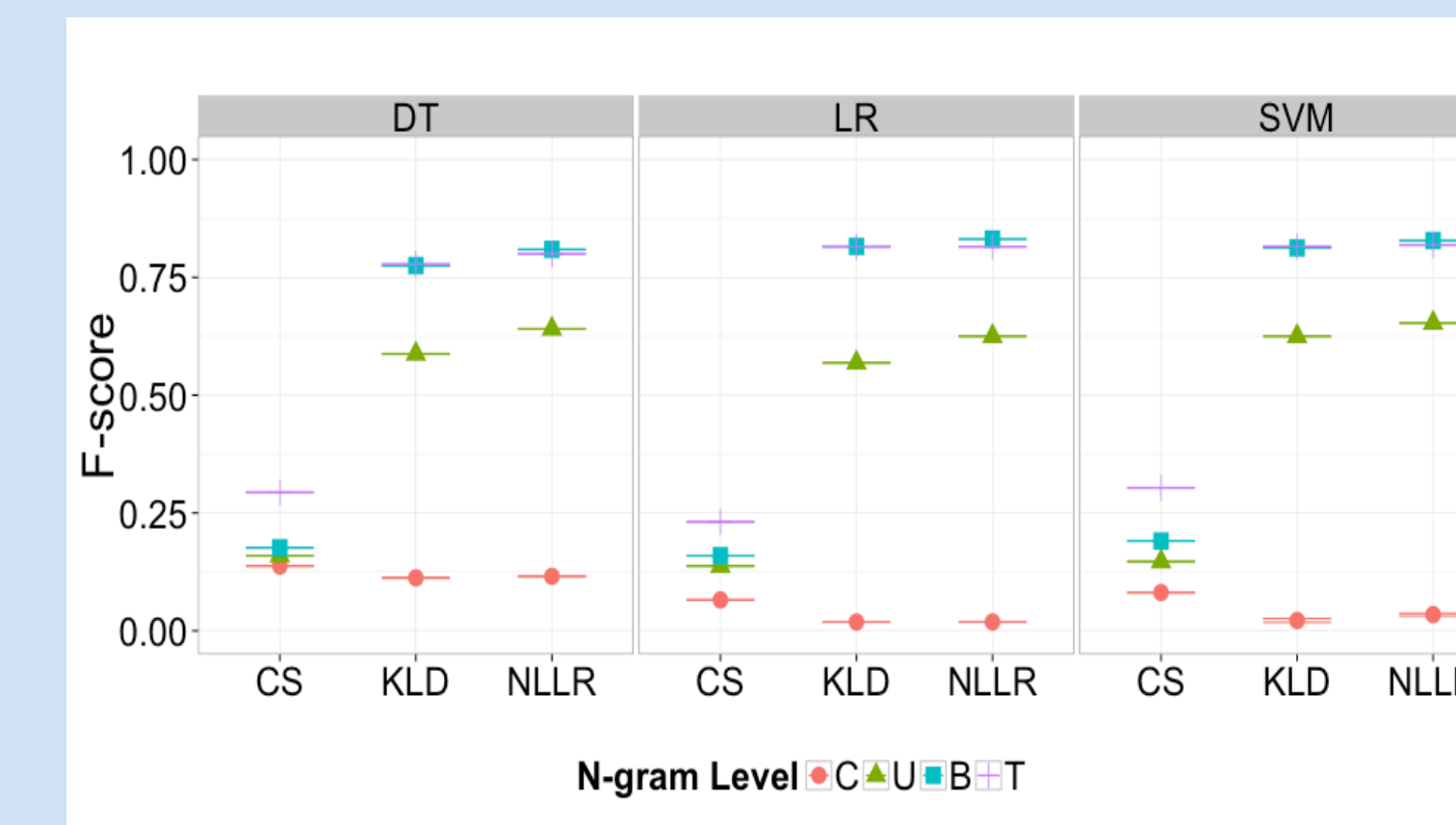
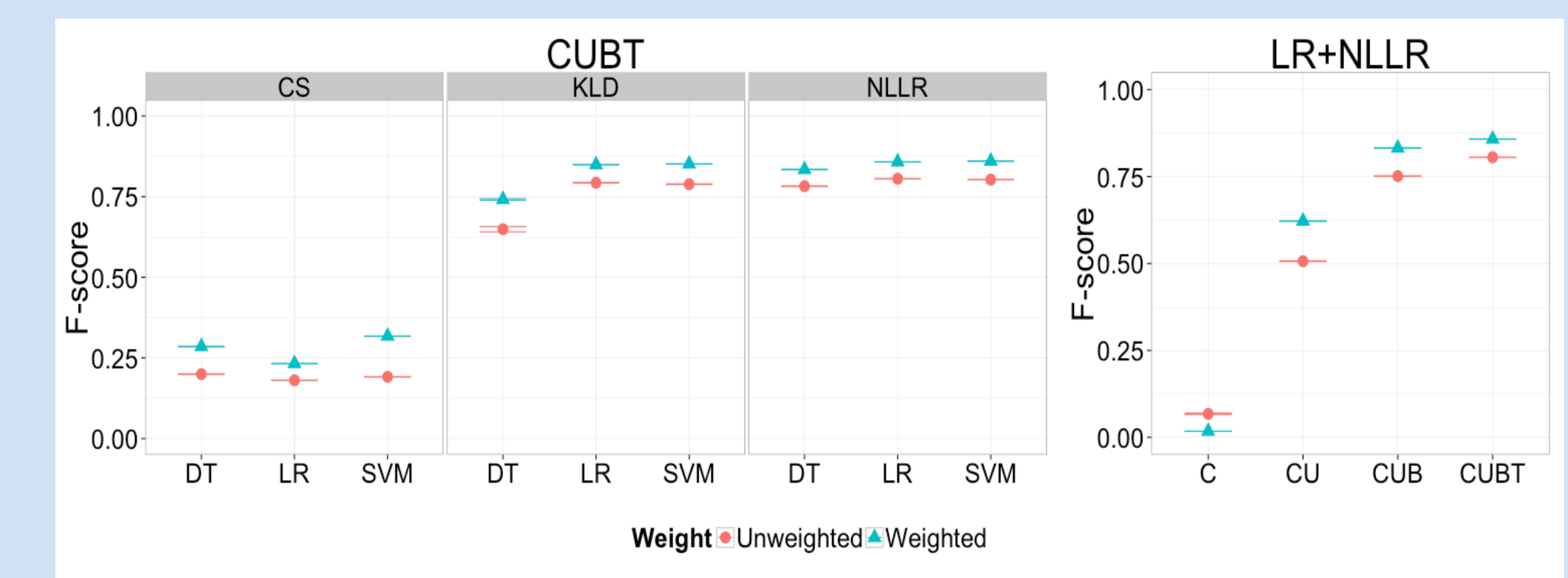


12 Chronons (bins) of books

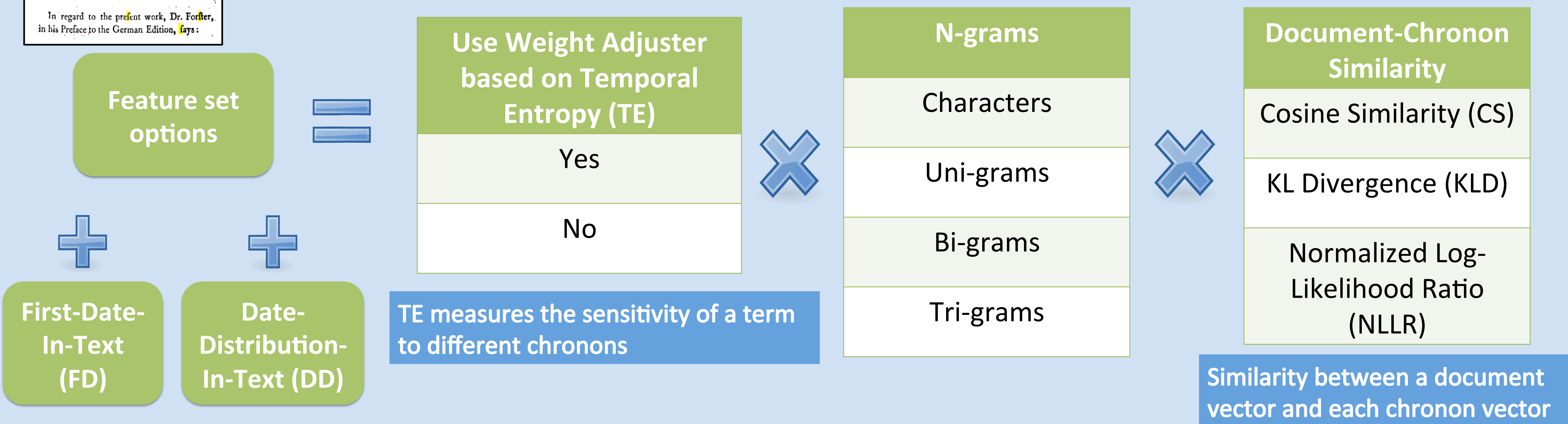
Experiment Design & Results

- ✓ Compare between classifiers
- ✓ Compare between n-grams
- ✓ Compare between similarities
- ✓ Compare between combinations of the above
- ✓ Observe baseline (FD feature only)

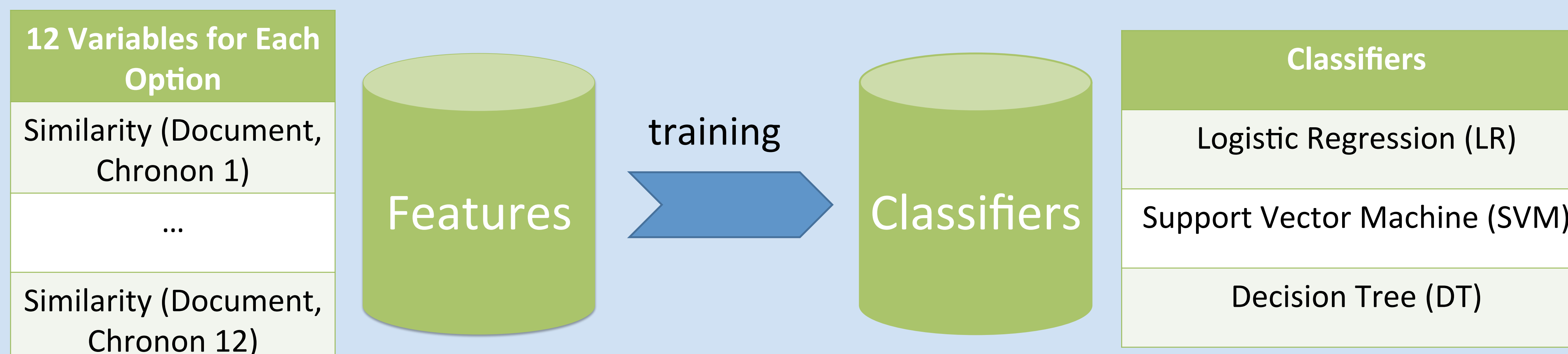
	BL	LR	DT	SVM
F-score	74.8	78.2	78.7	77.0
Precision	79.1	79.9	78.8	79.3
Recall	73.2	77.7	78.7	76.4



Feature Sets



Machine Learning



Takeaways

- Higher-order n-grams like bigrams and trigrams are more effective than unigrams
- Proposed bag-of-character to capture OCR errors as a feature
- Temporal entropy (TE) is an effective term weighting strategy
- The methodology works reasonably well on HT digital library data
- Can generalize to other large-scale digital library metadata records

Acknowledgement to HathiTrust Research Center for the NGPD data set and metadata.