

Fusing Data with Correlations

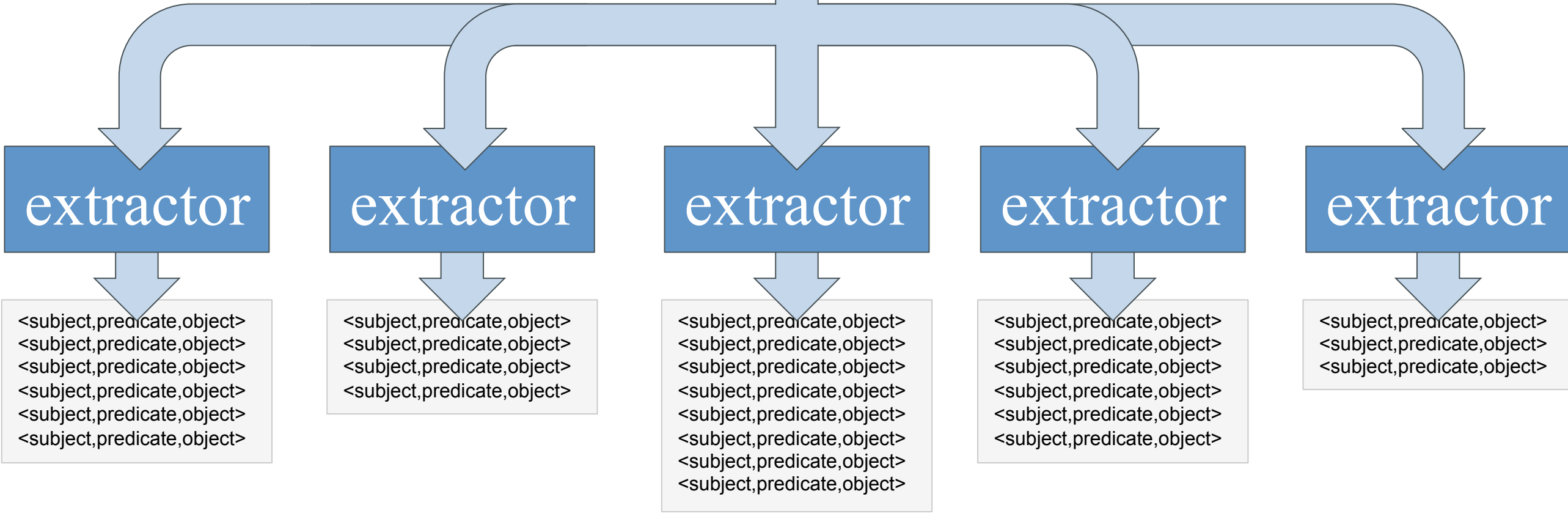
Ravali Pochampally
UMass Amherst

Anish Das Sarma
Troo.ly

Luna Dong
Google

Alexandra Meliou
UMass Amherst

Divesh Srivastava
AT&T Research



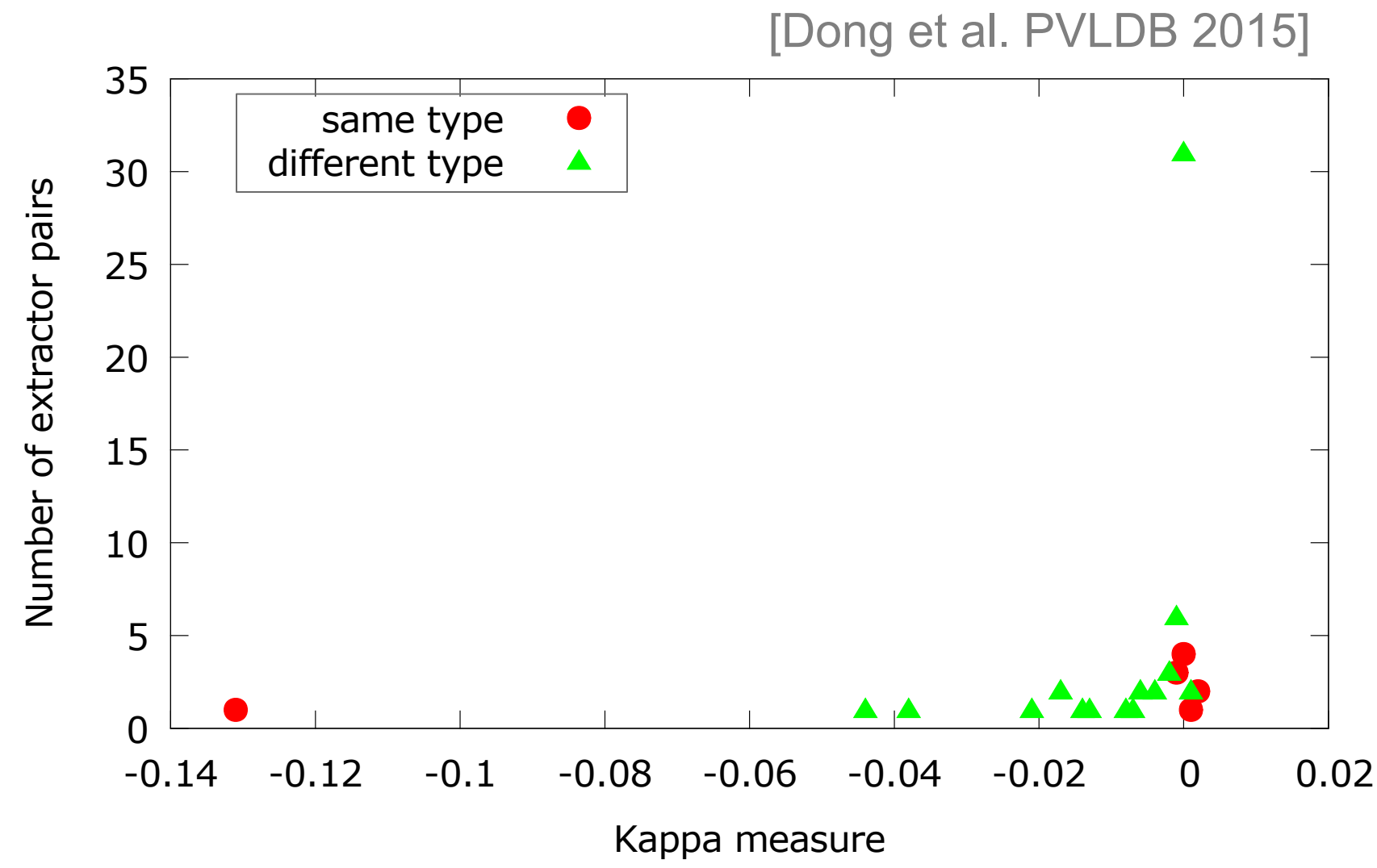
Knowledge base

ID	KnowledgeTriple	S ₁	S ₂	S ₃	S ₄	S ₅
t ₁	{Obama, profession, president}	✓	✓		✓	✓
t ₂	{Obama, died, 1982}	✓	✓			
t ₃	{Obama, profession, lawyer}			✓		
t ₄	{Obama, religion, Christian}		✓	✓	✓	✓
t ₅	{Obama, age, 50}		✓	✓		
t ₆	{Obama, support, White Sox}	✓			✓	✓
t ₇	{Obama, spouse, Michelle}	✓	✓	✓		
t ₈	{Obama, administered by, John G. Roberts}	✓	✓		✓	✓
t ₉	{Obama, surgical operation, 05/01/2011}	✓	✓		✓	✓
t ₁₀	{Obama, profession, community organizer}	✓		✓	✓	✓

anti-correlated

correlated

Correlations in web extraction

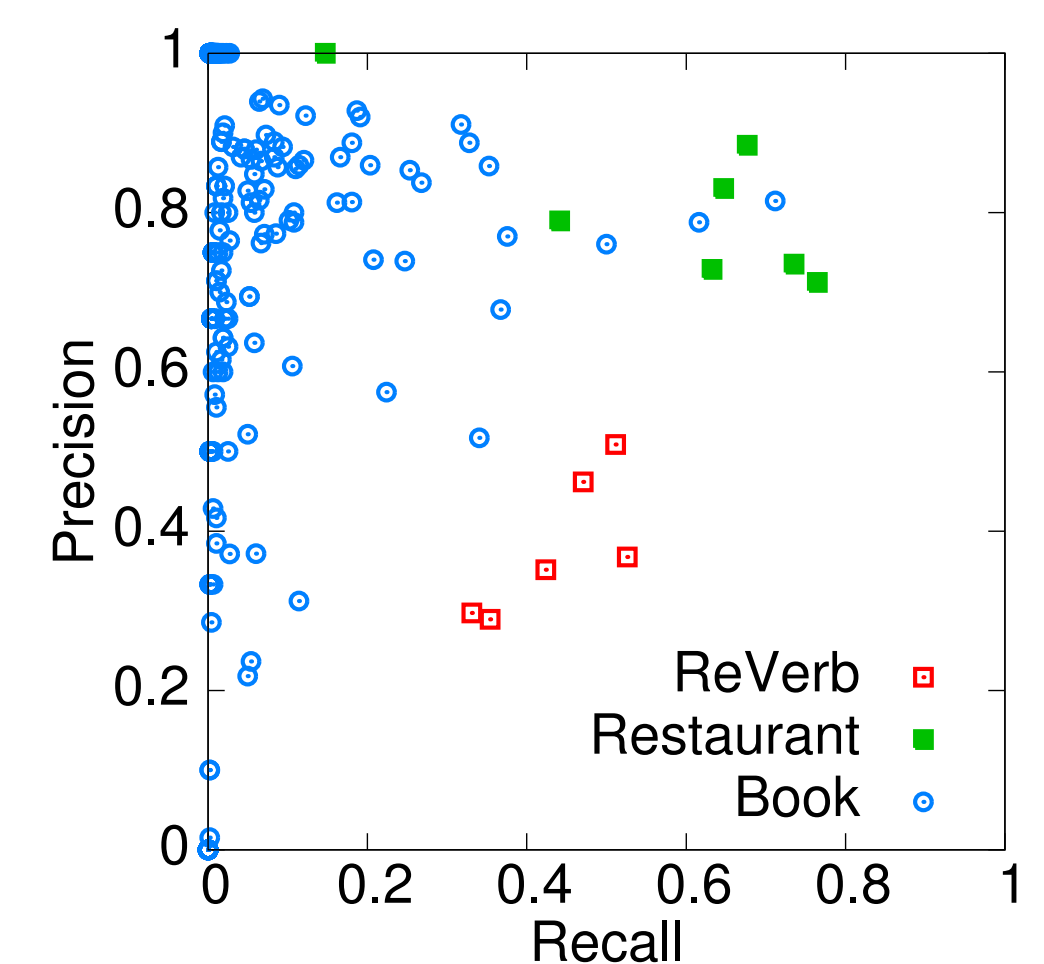


Correlations are richer than copying relationships

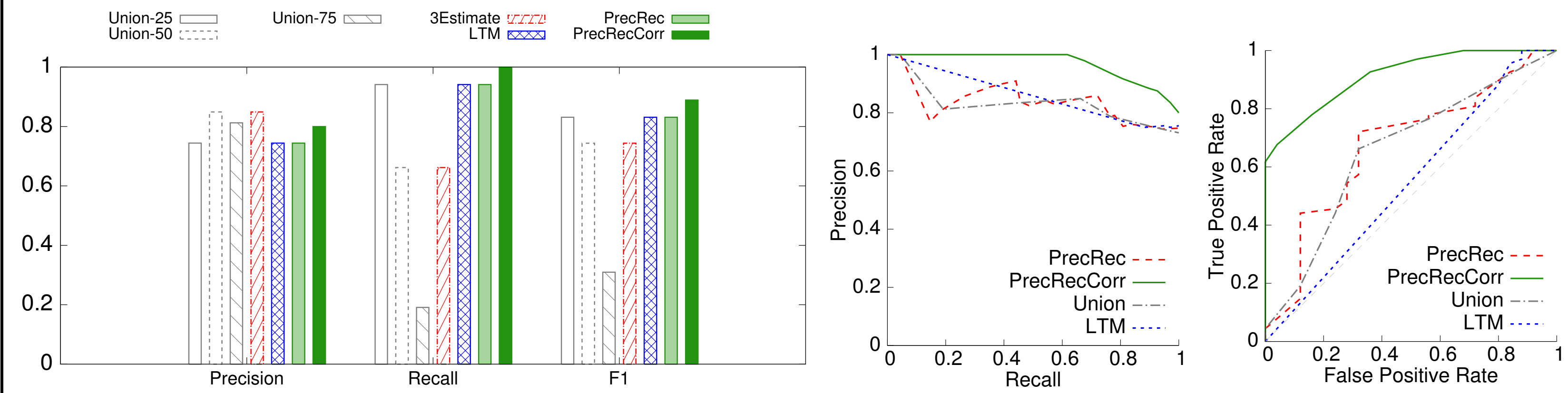
There is an exponential number of correlation parameters, but we provide a scalable solution

Real-world datasets

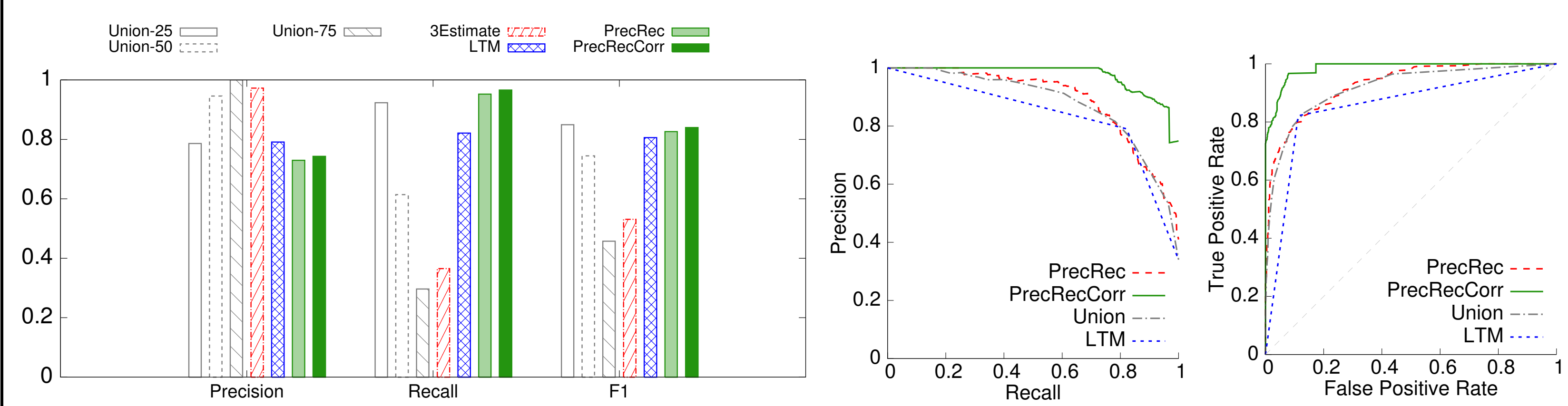
- Union-K
– a triple is true if at least K sources provide it
- 3-Estimate [Galland et al. WSDM 2010]
– iteratively computes trustworthiness
- LTM [Zhao et al. VLDB 2012]
– uses graphical models and Gibbs sampling
- PrecRec / PrecRecCorr **our methods**



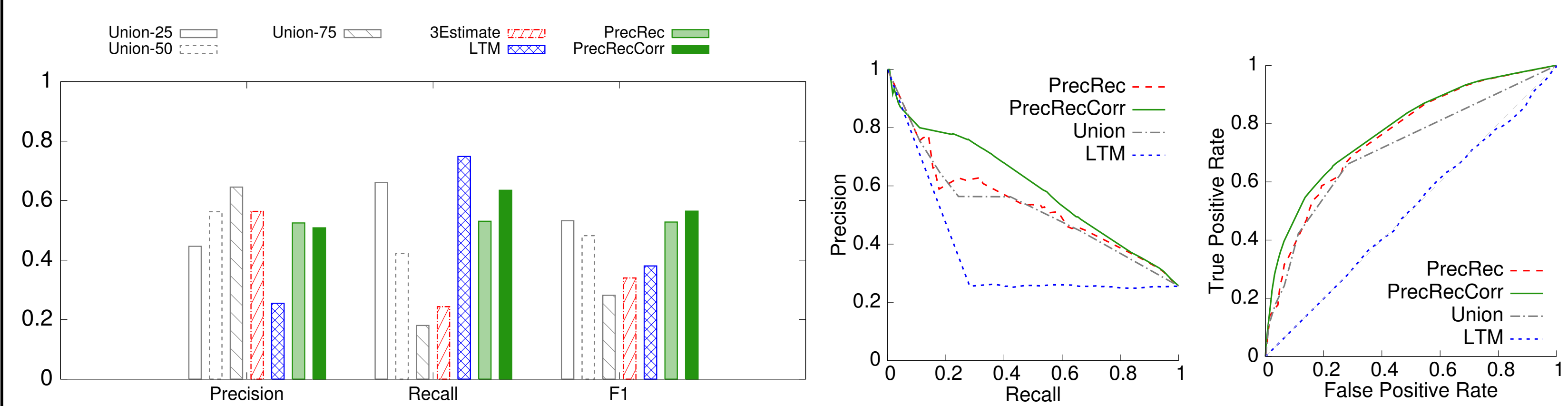
Restaurant: 7 sources, 93 triples



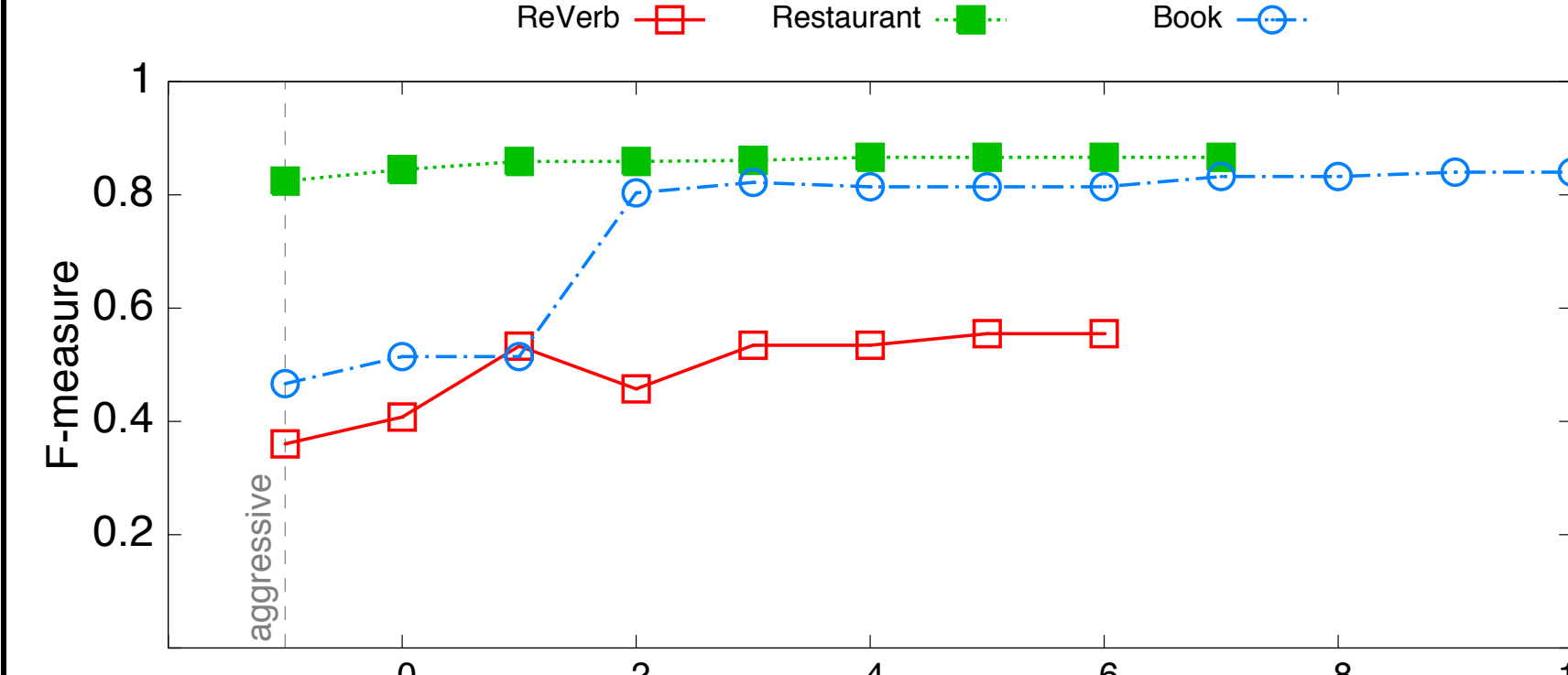
Book: 879 sources, 225 triples



ReVerb: 6 extractors, 2407 triples



elastic approximation

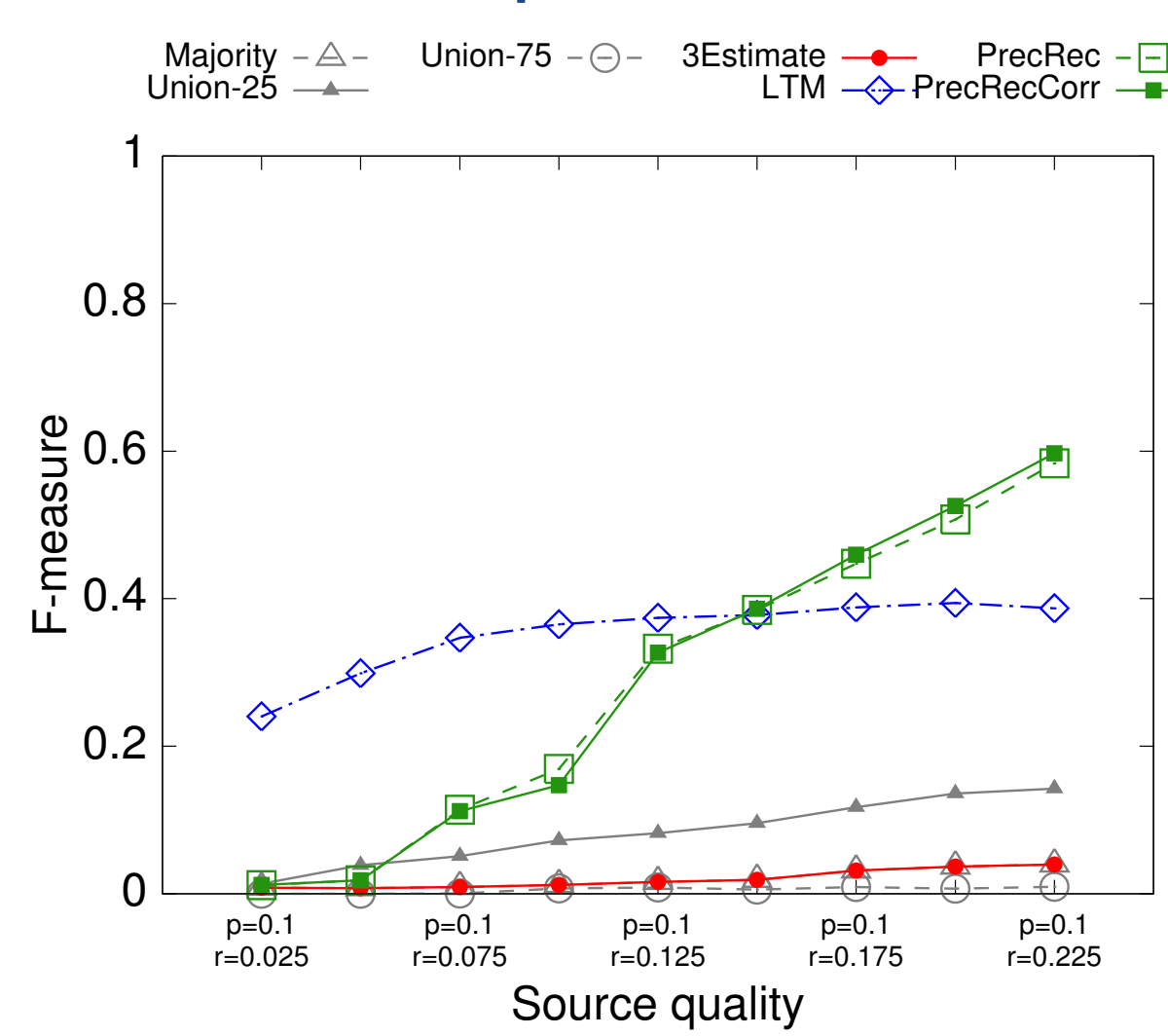


runtime comparison

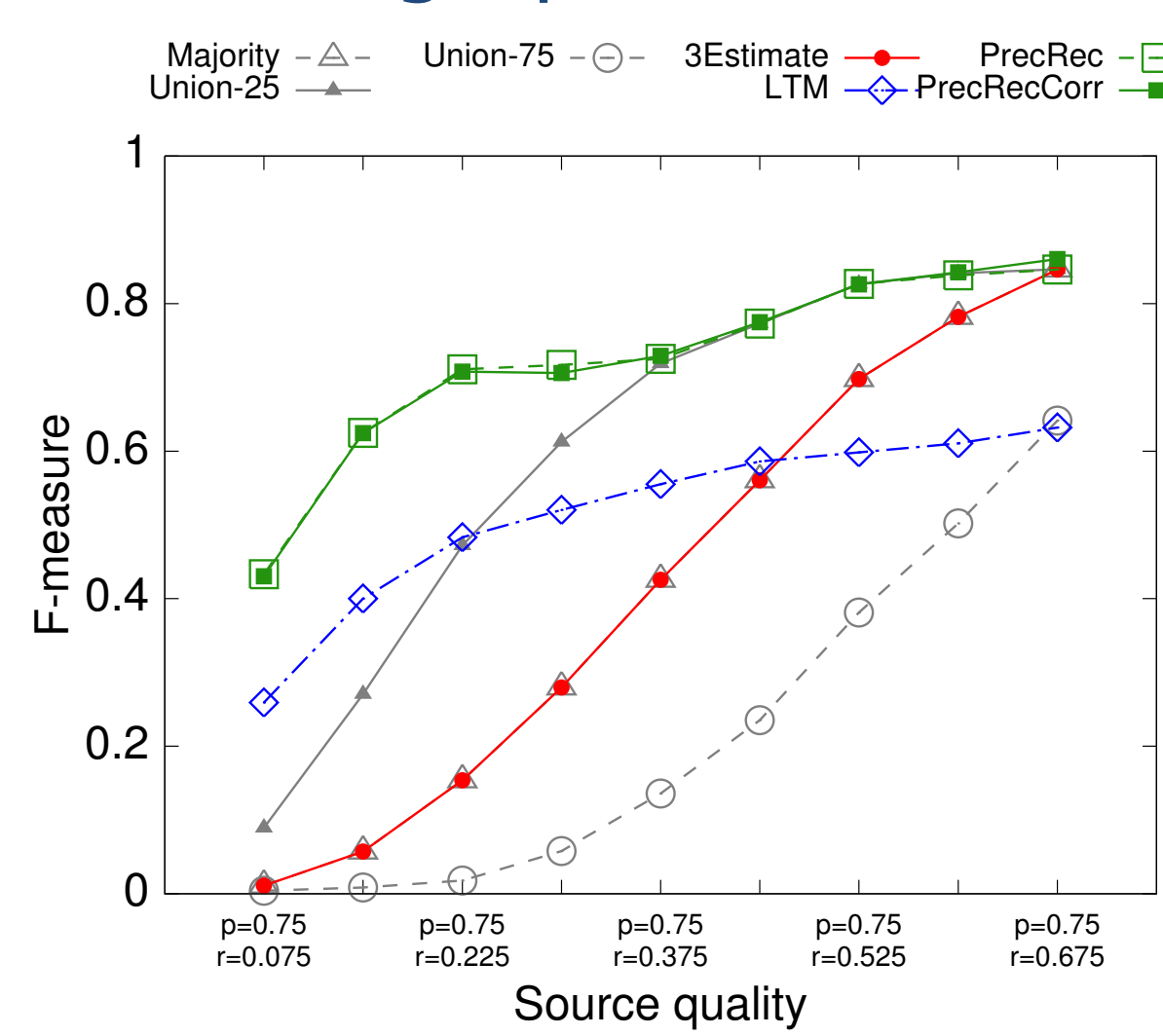
time(sec)	REVERB	RESTAURANT	BOOK
UNION-25	0.39	0.56	3.86
UNION-50	0.14	0.32	3.71
UNION-75	0.11	0.35	3.00
3-ESTIMATE	0.7	0.06	39
LTM (10 iter)	49	5.3	3791
PRECREC	2.6	0.3	35
PRECRECCORR	124	5.4	6786
PRECRECCORR-LVL3	79	2.25	2452

Synthetic datasets

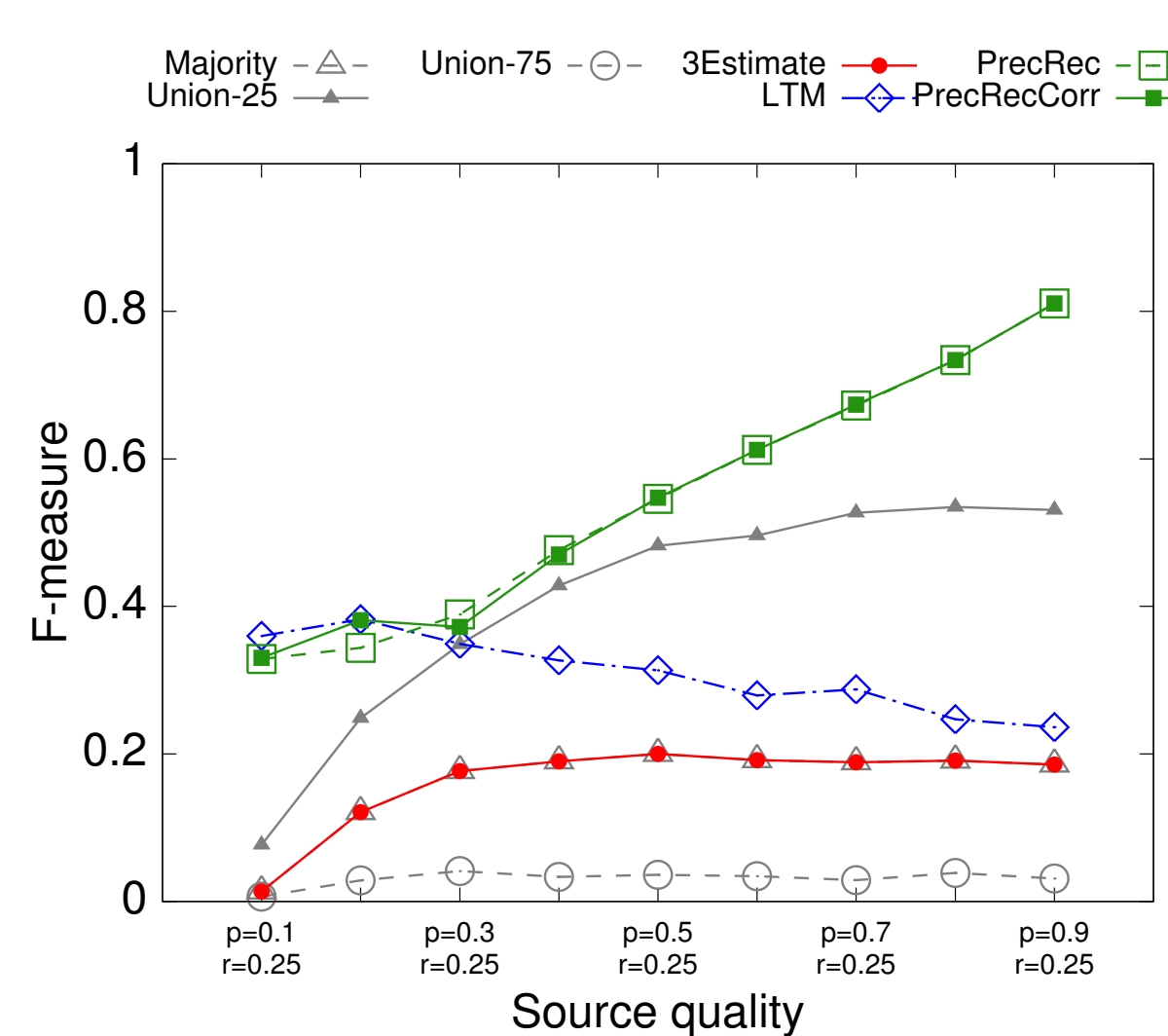
low precision



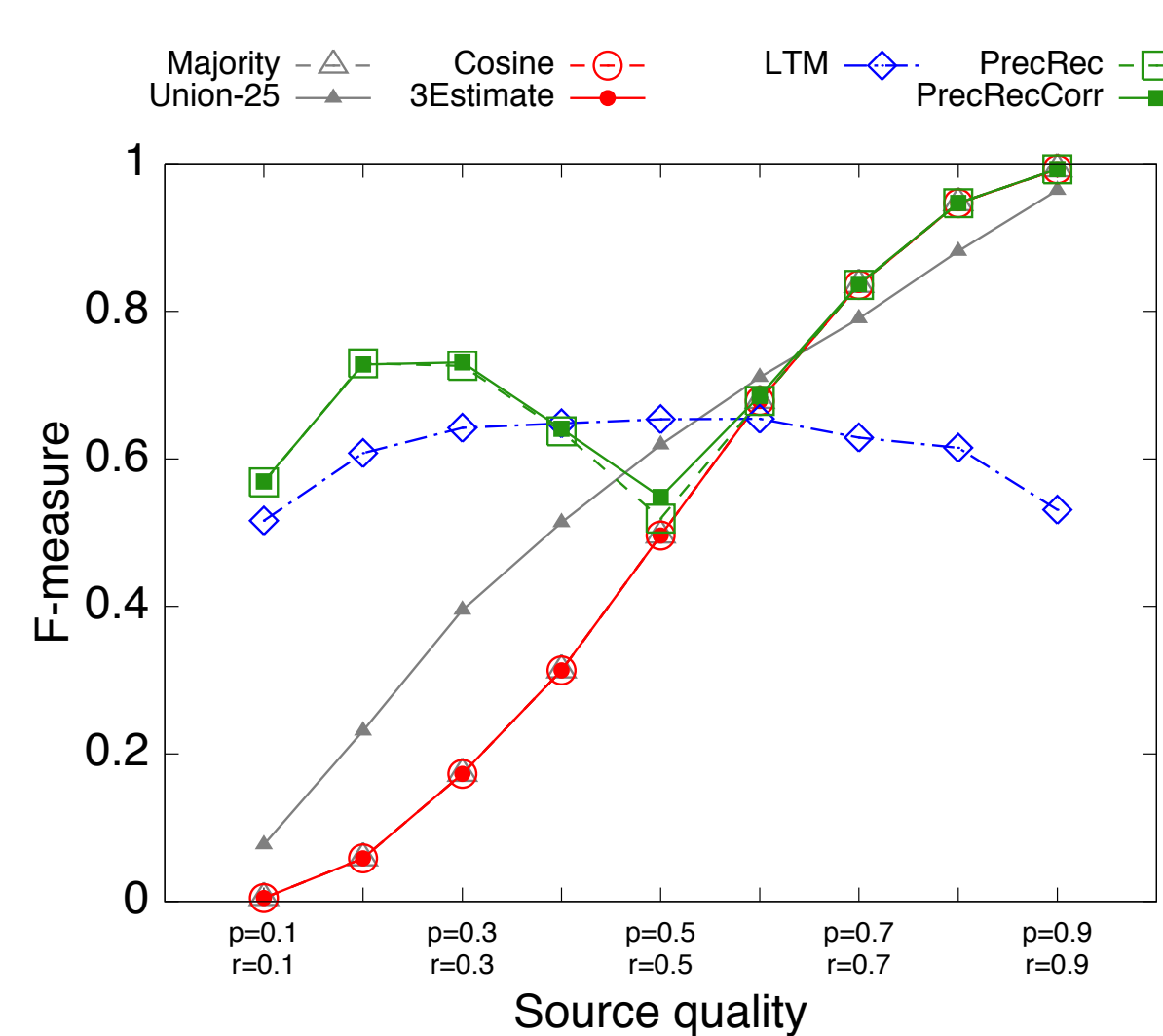
high precision



low recall



bad to good quality



correlation

