

Configuring Topic Models for Software Engineering Tasks in TraceLab

Bogdan Dit

Annibale Panichella

Evan Moritz

Rocco Oliveto

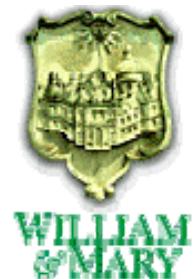
Massimiliano Di Penta

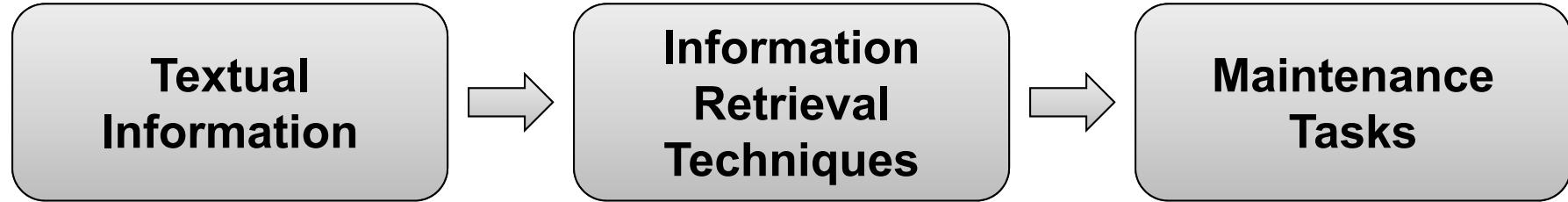
Denys Poshyvanyk

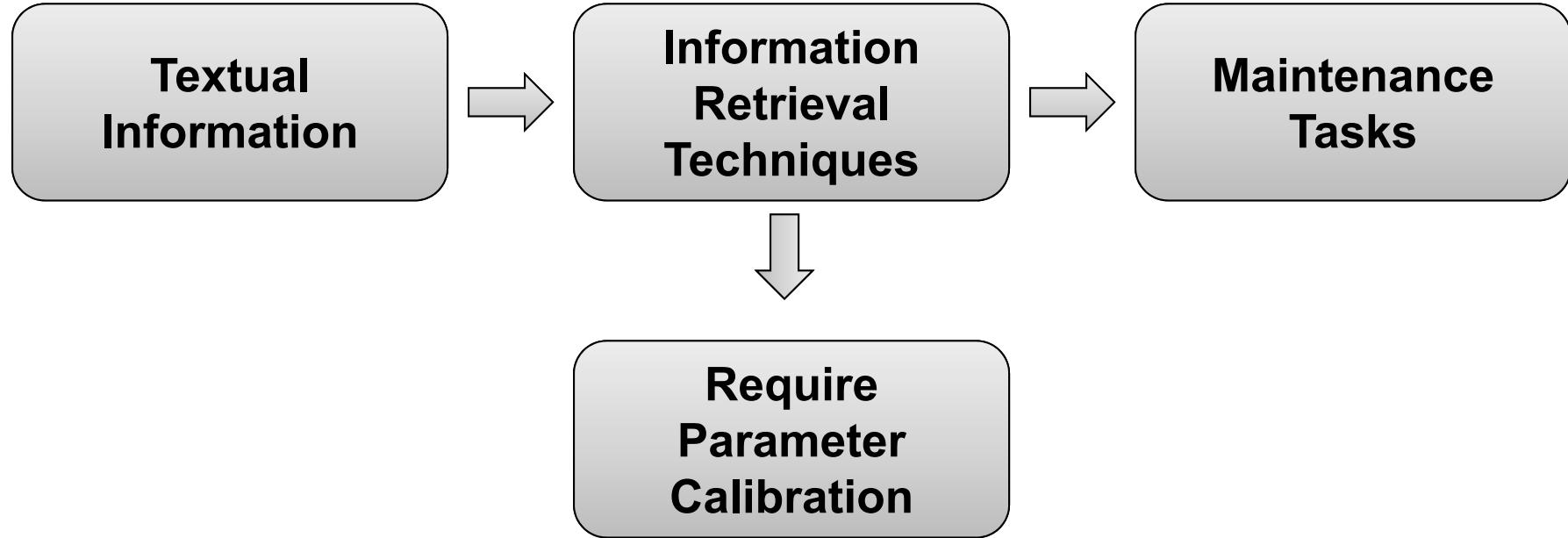
Andrea De Lucia

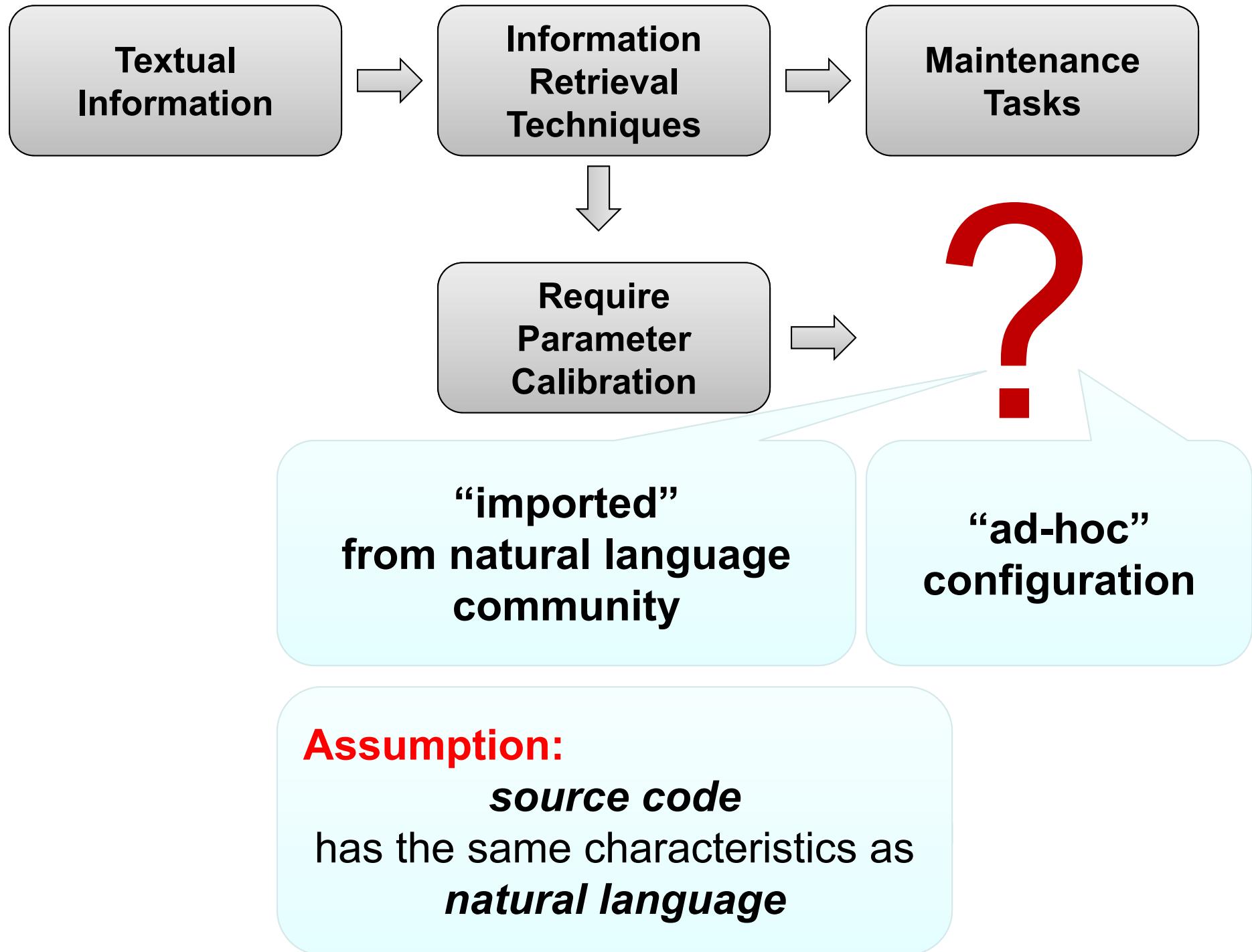


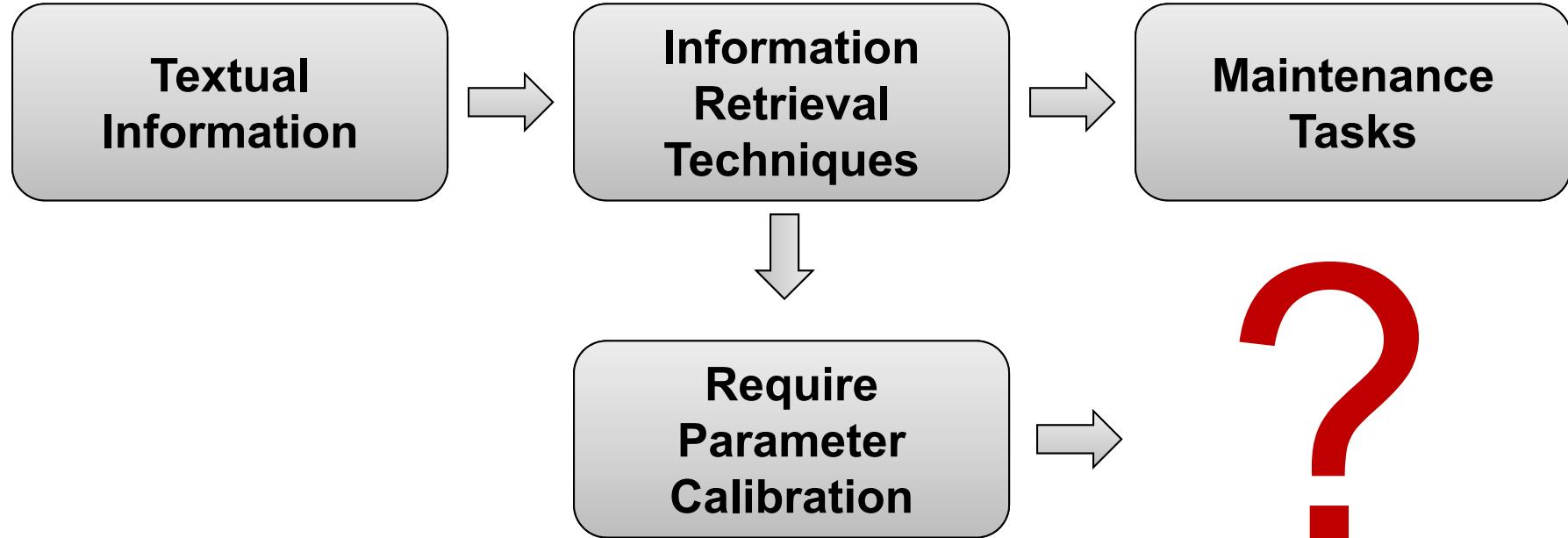
TEFSE'13
San Francisco, CA

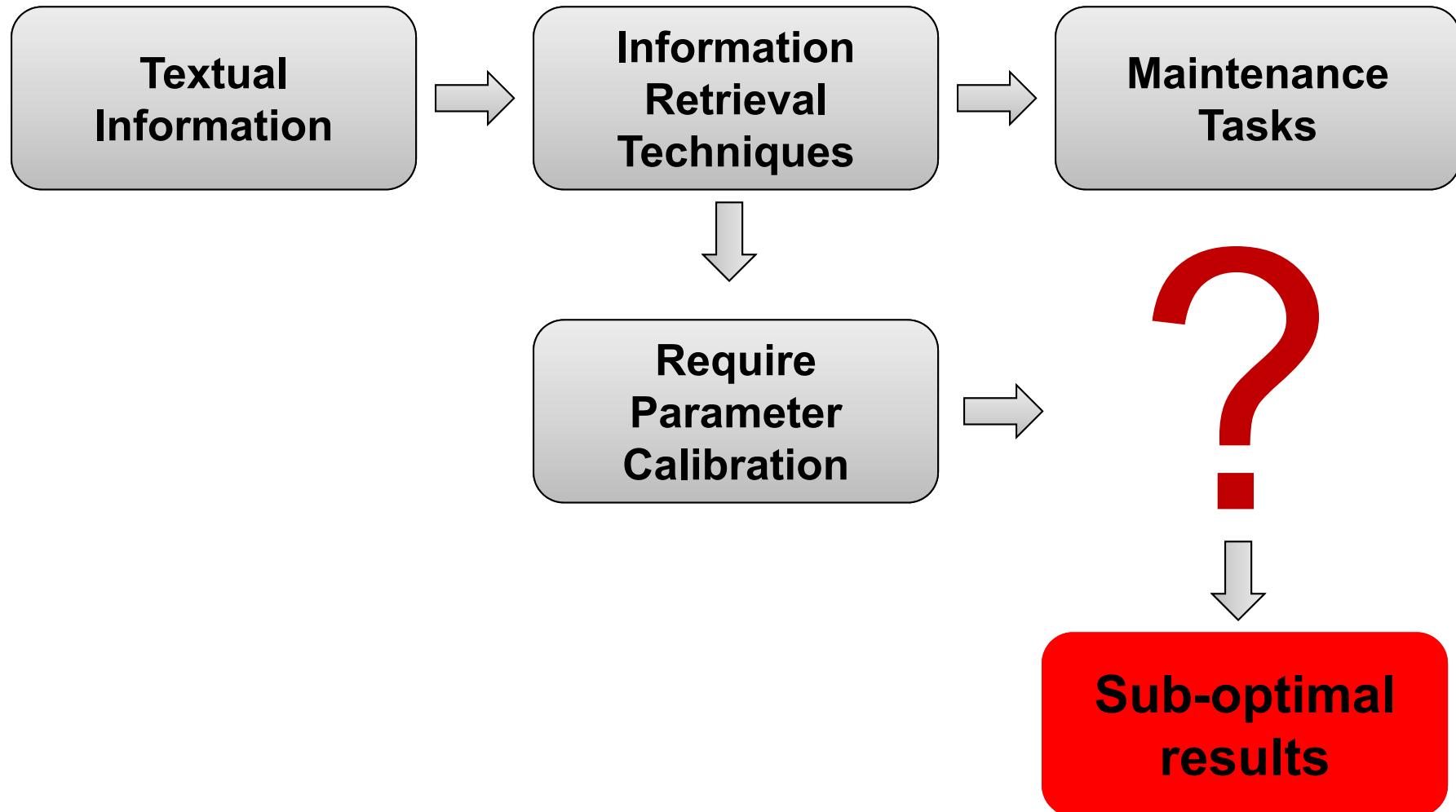


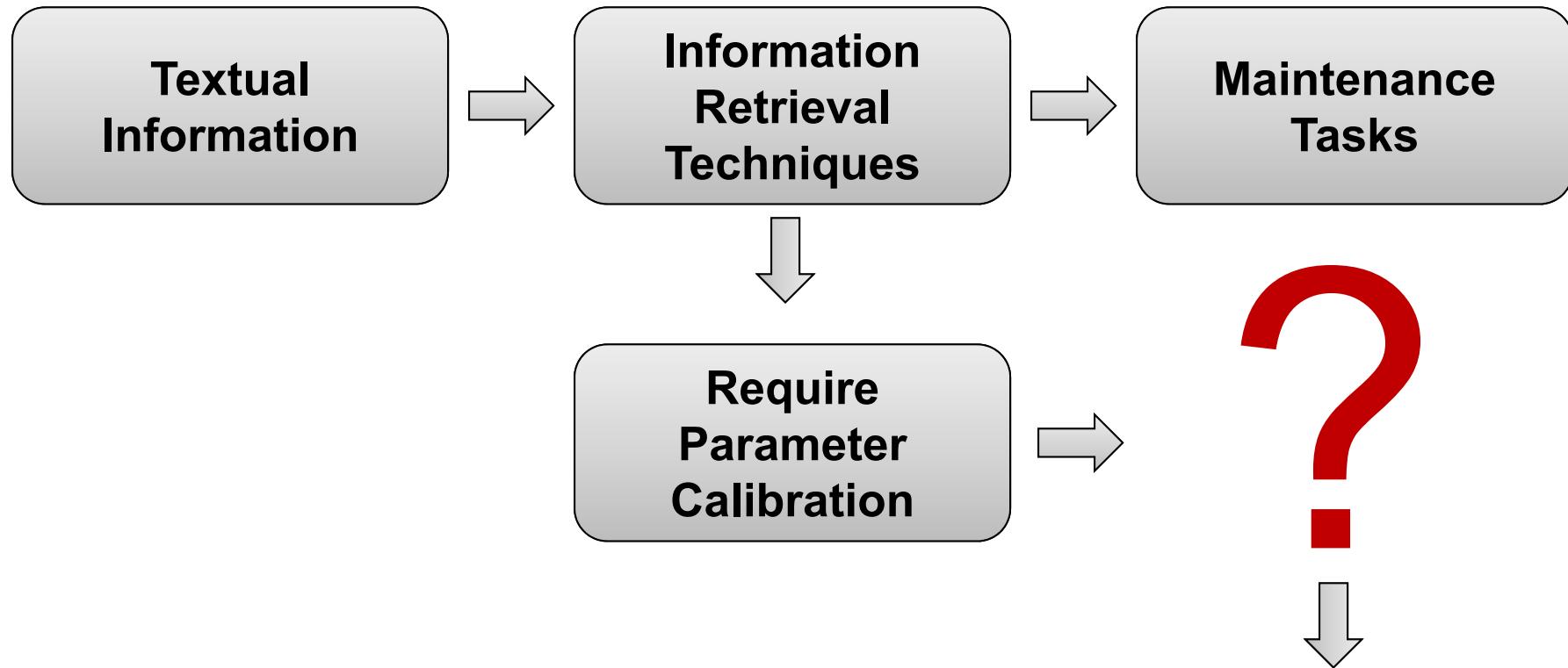




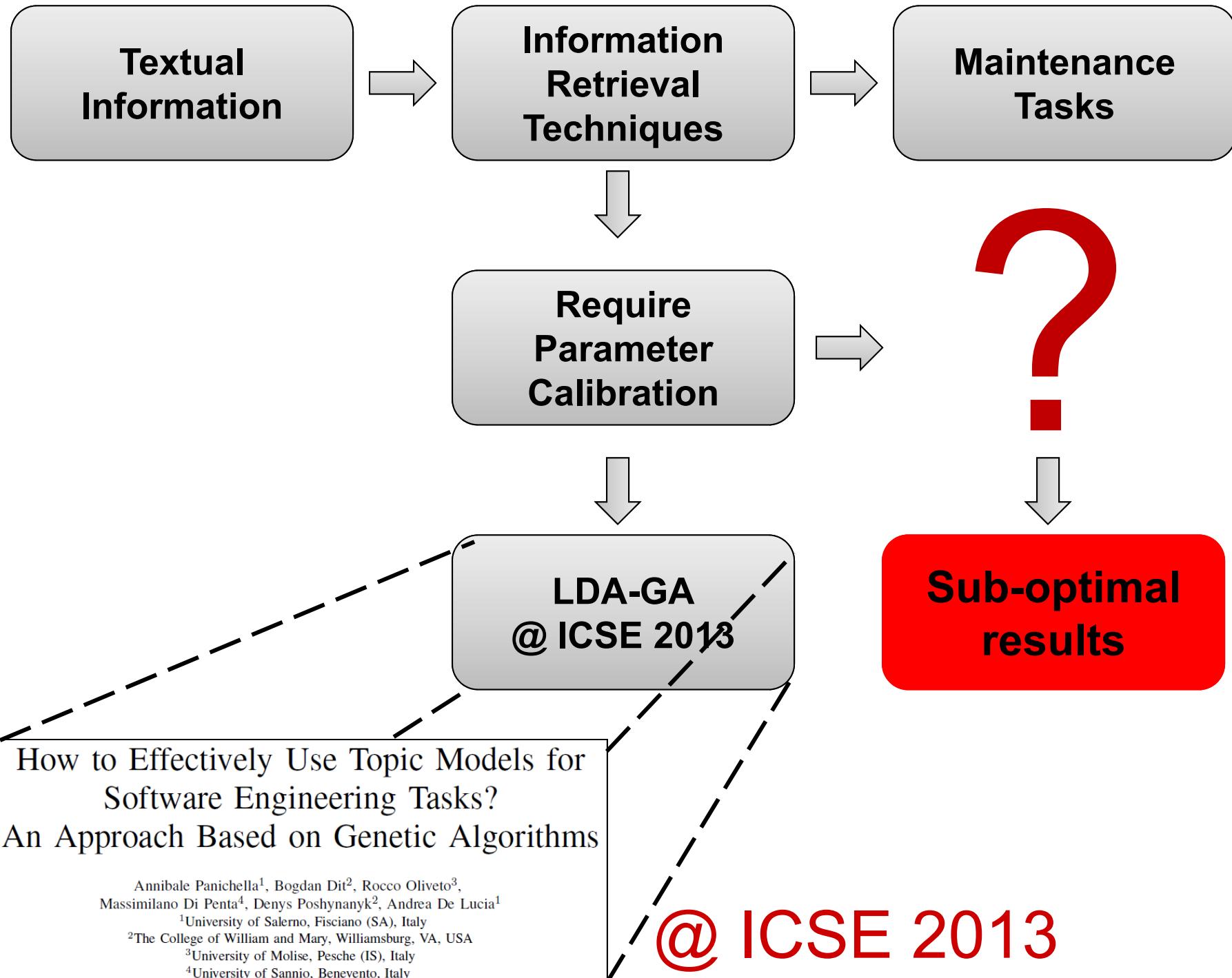


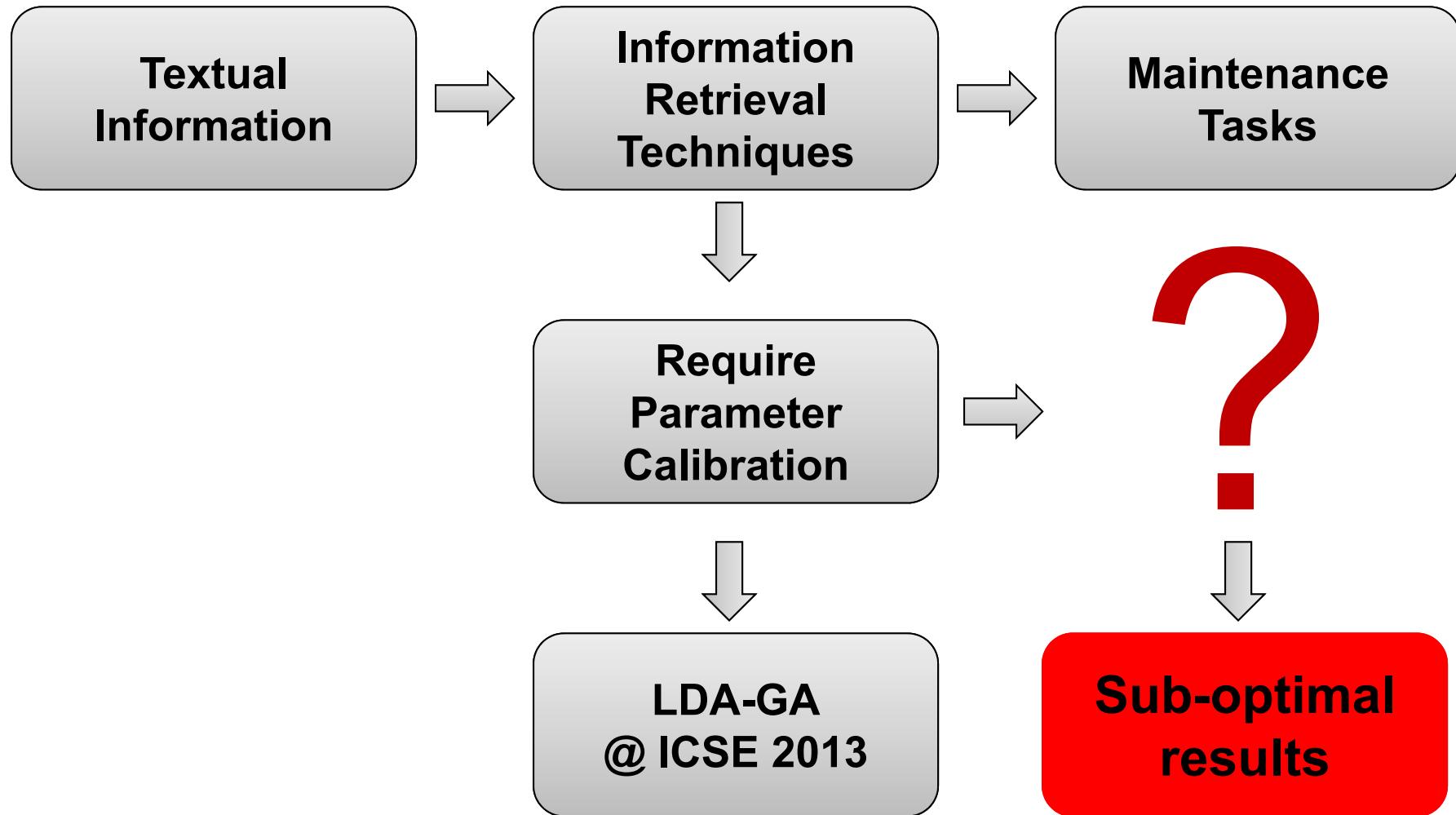


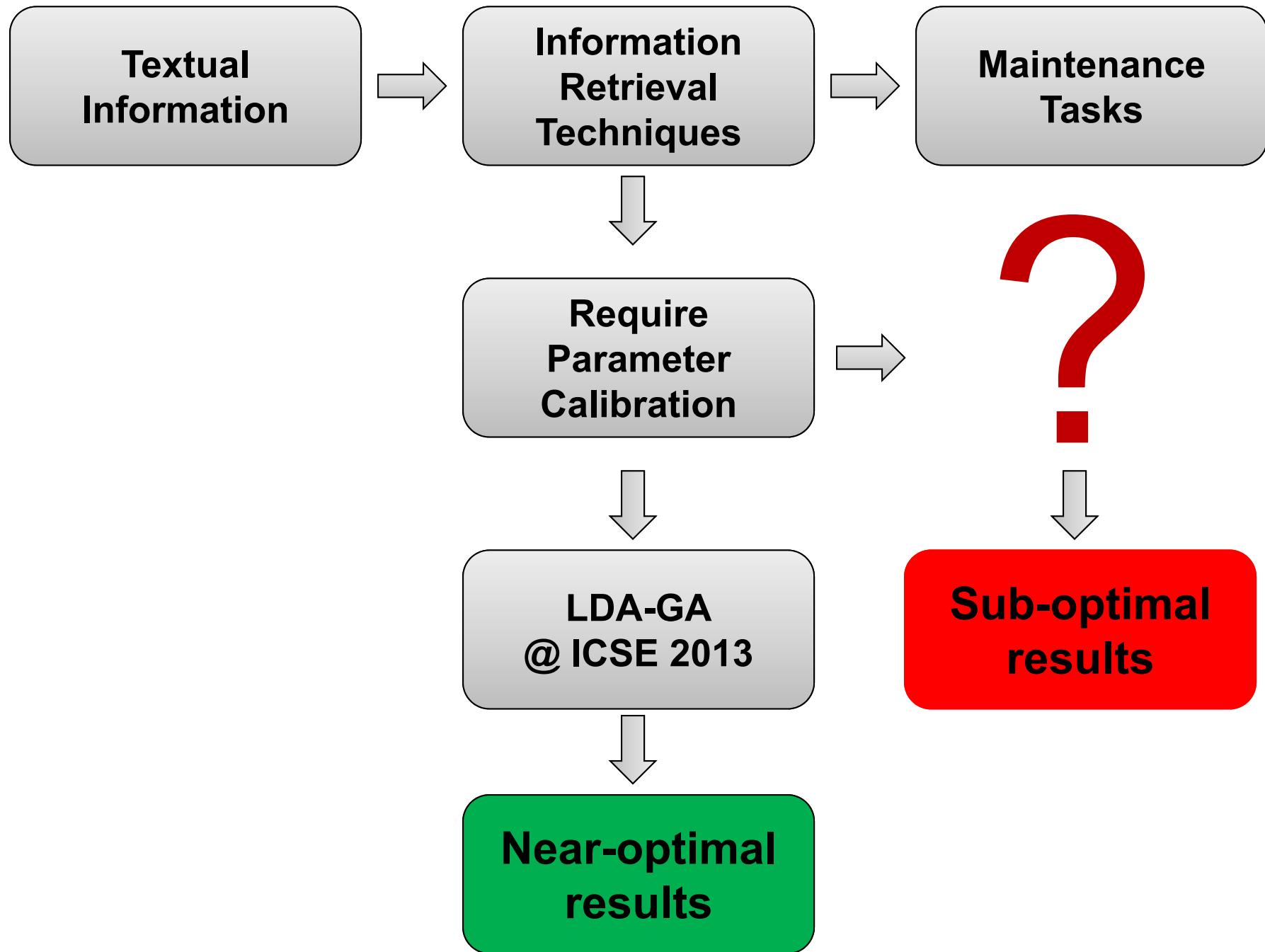


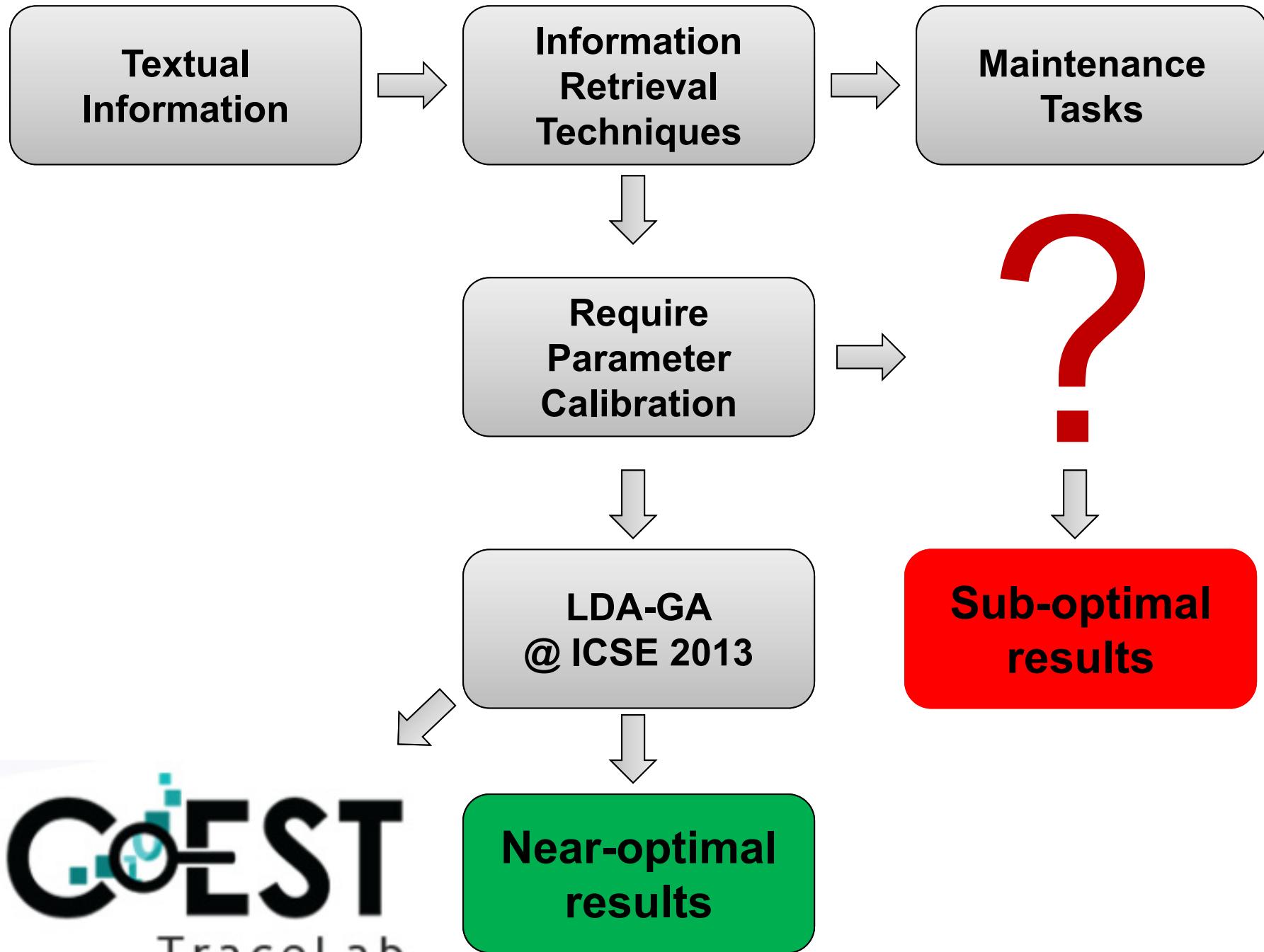


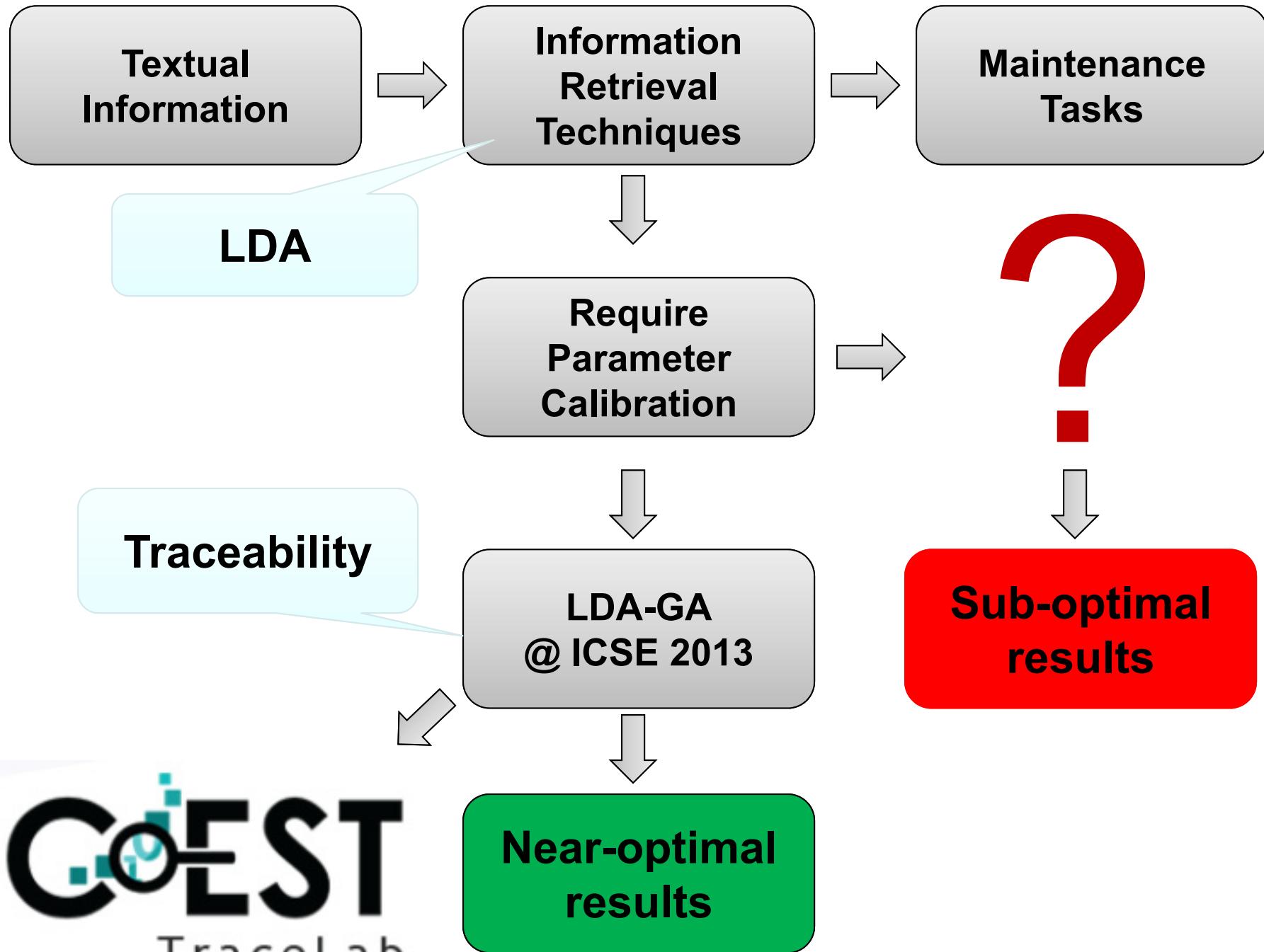
[Hindle et al. @ ICSE'12]:
source code
is more **regular** and **predictable** than
natural language











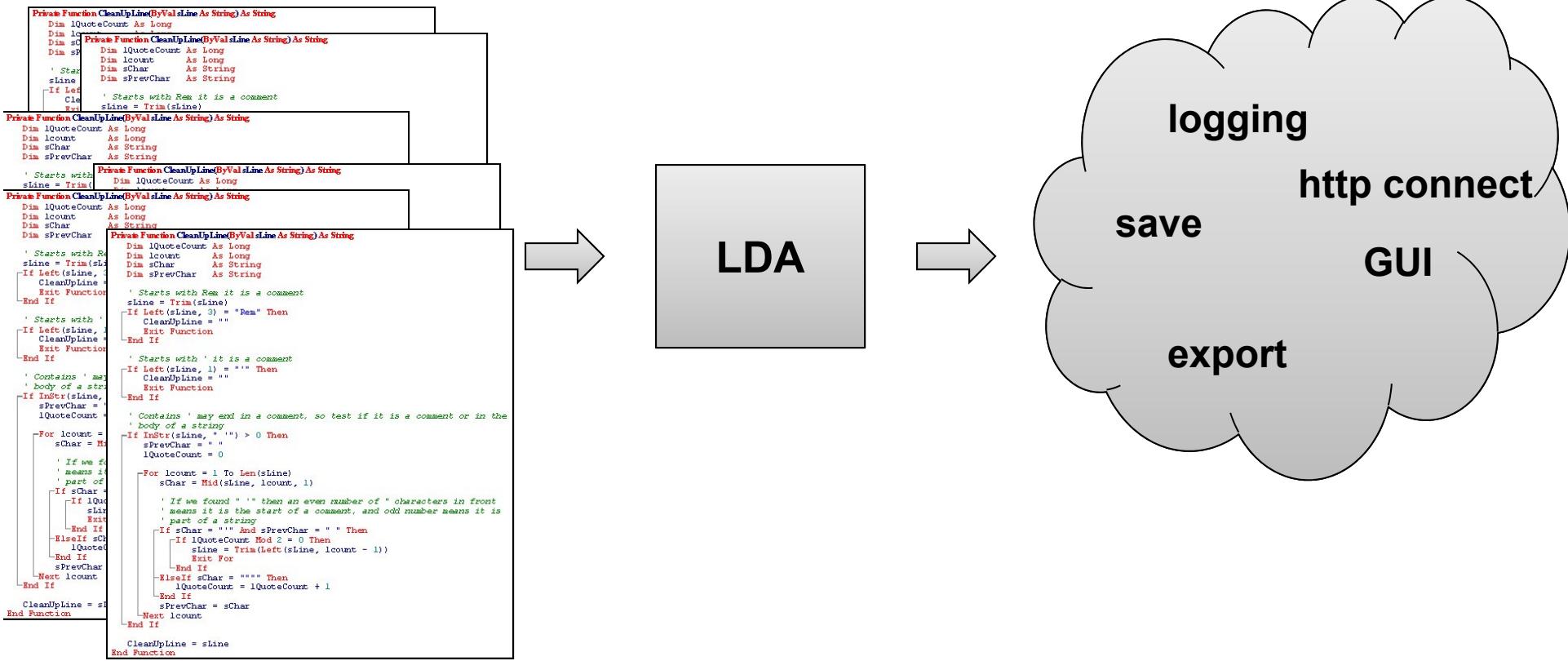
What is LDA?

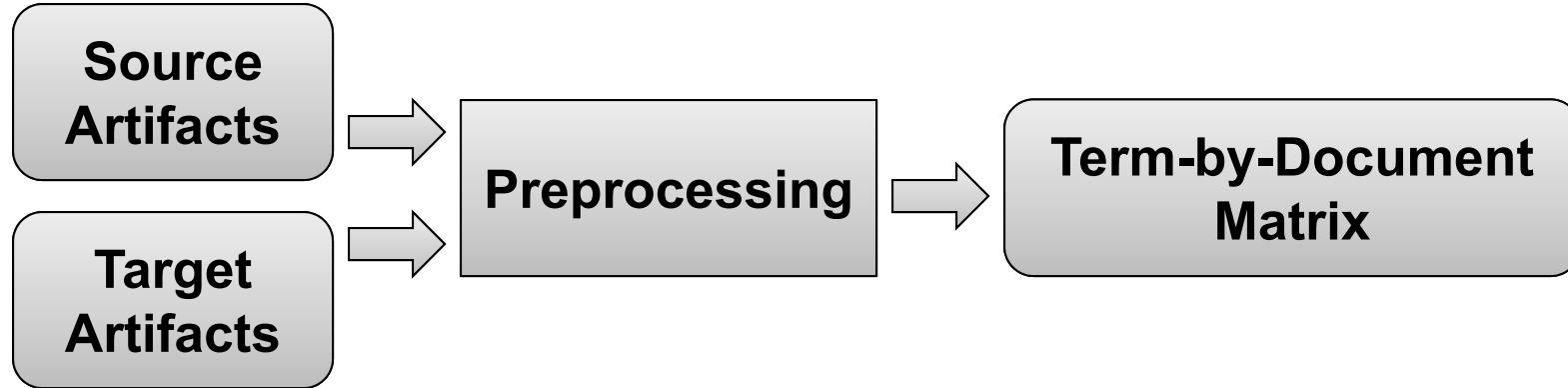
Latent Dirichlet Allocation (LDA)

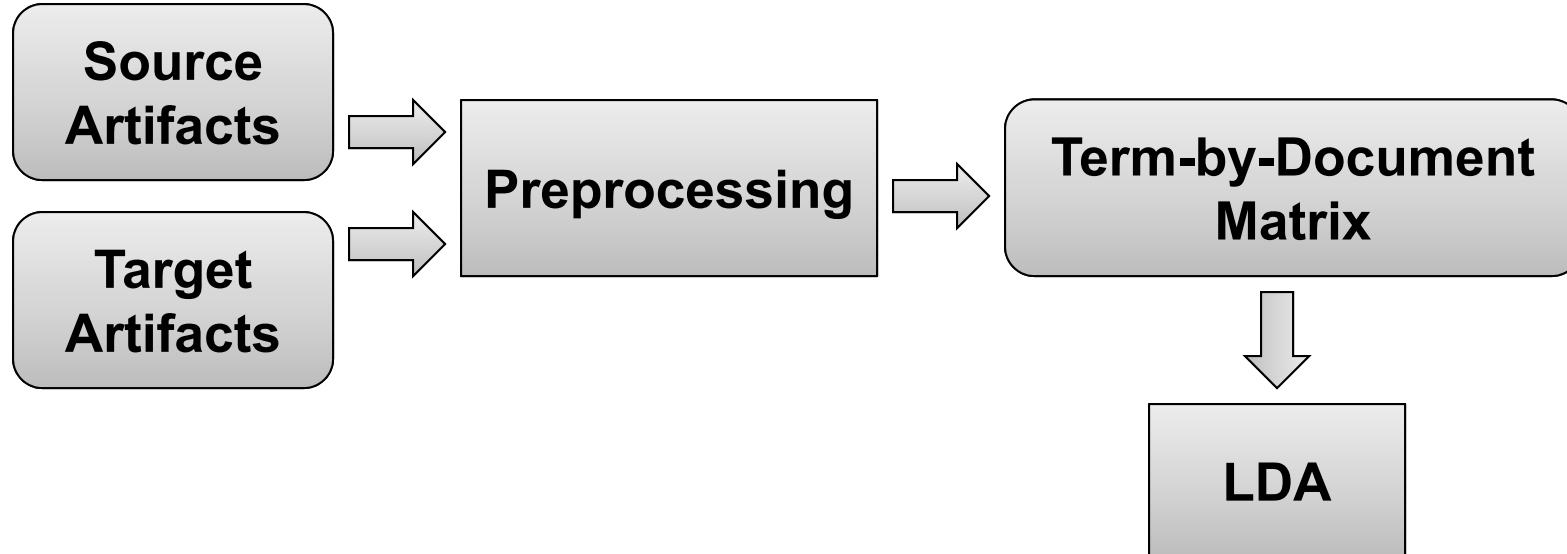
- Topic model that generates the distribution of latent topics from textual documents

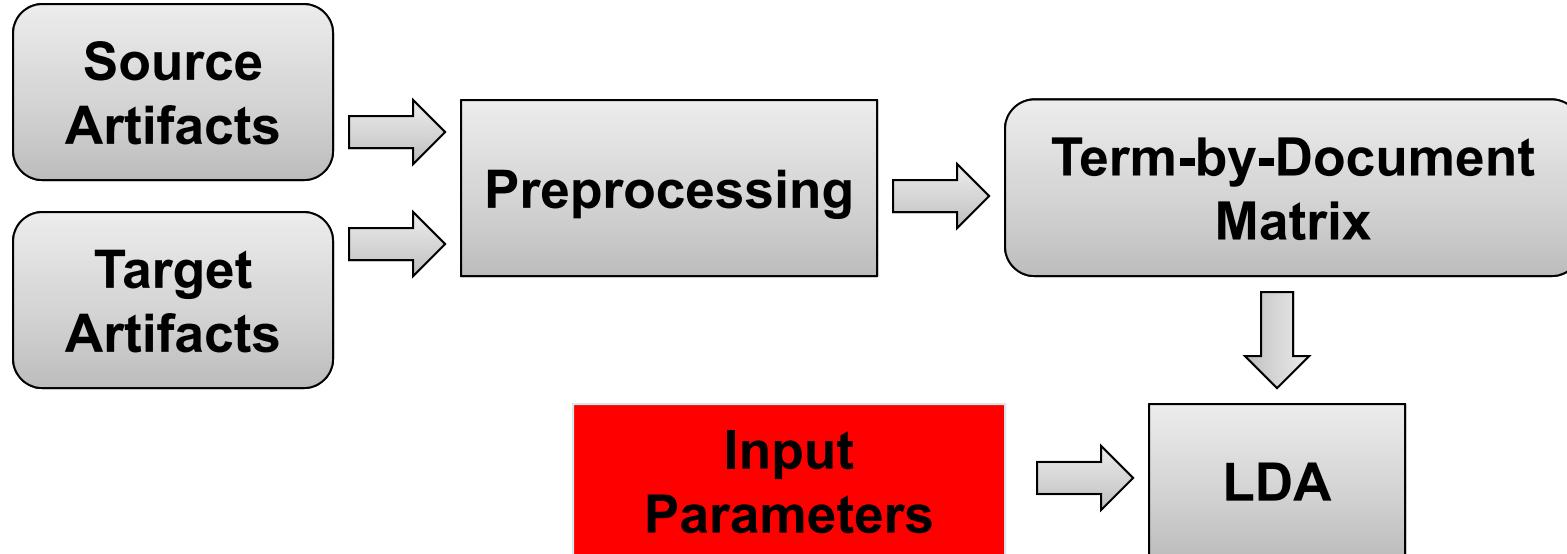
Latent Dirichlet Allocation (LDA)

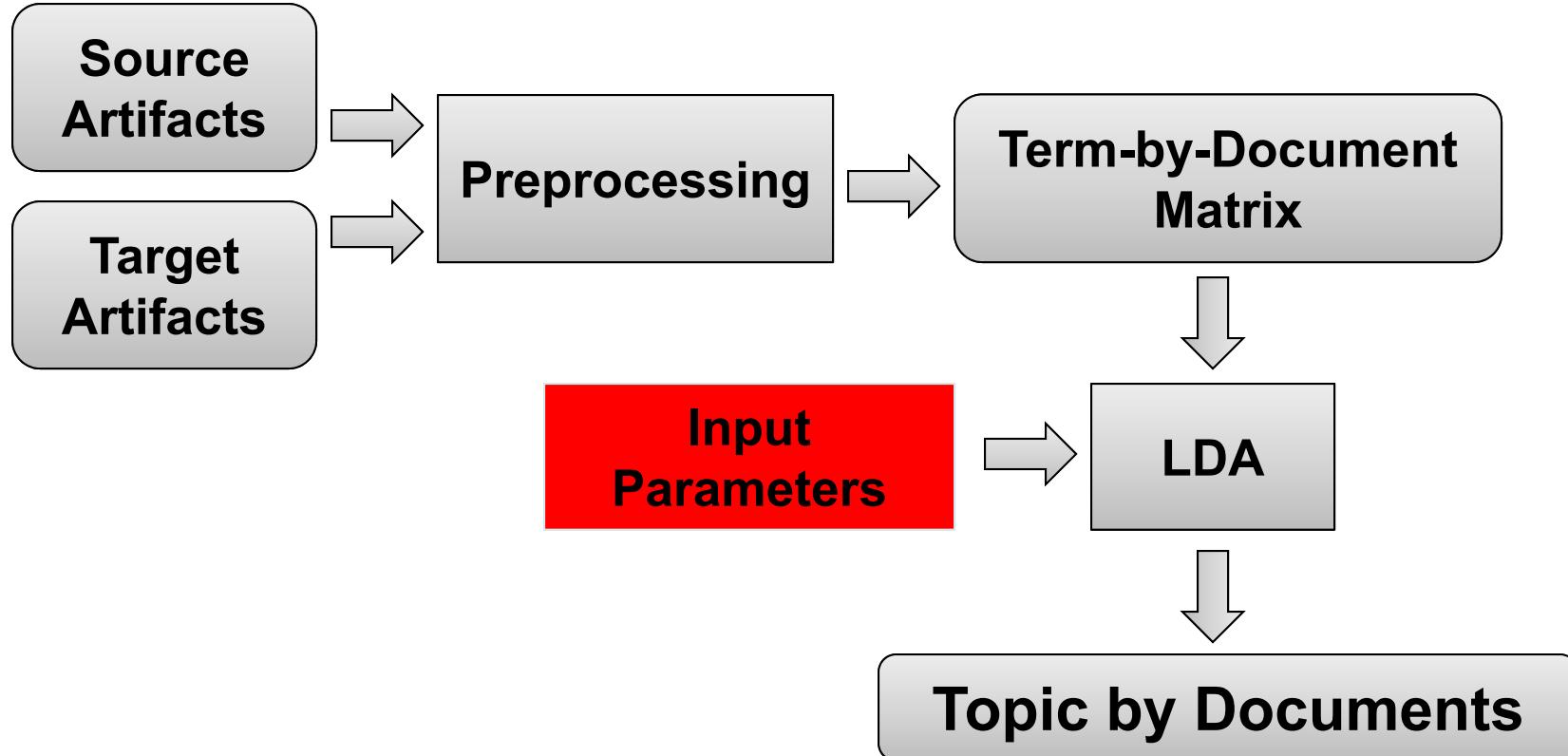
- Topic model that generates the distribution of latent topics from textual documents





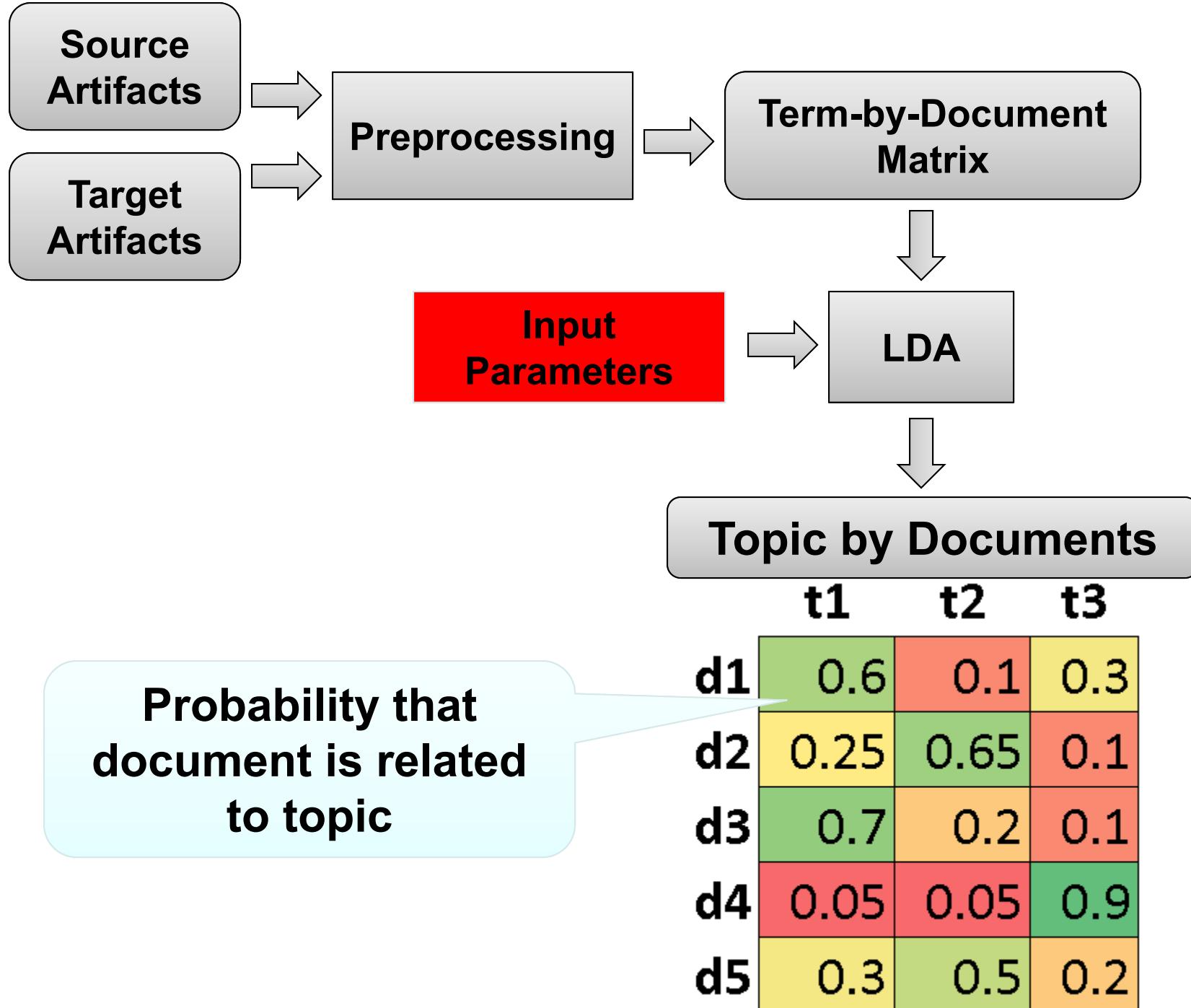


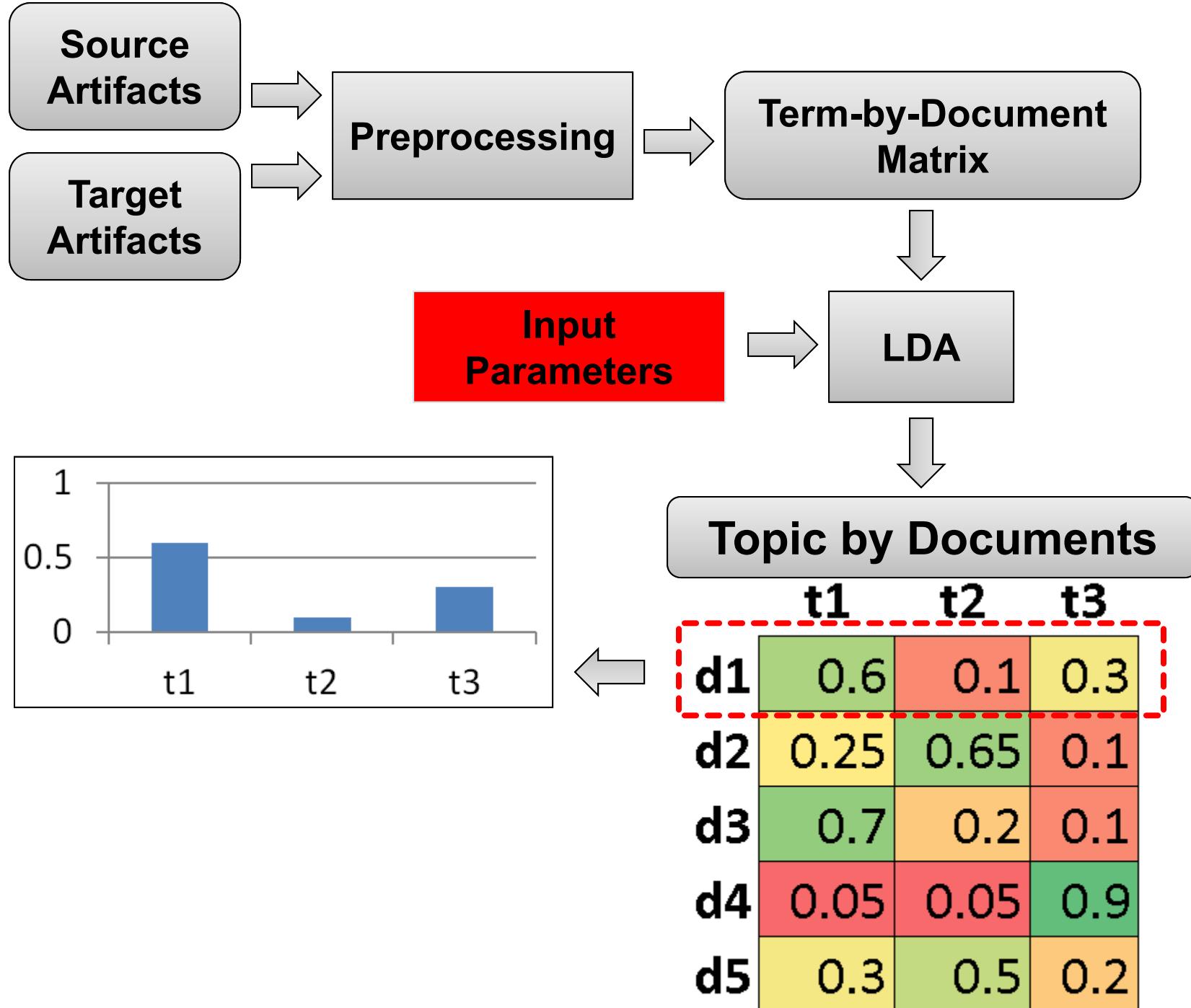


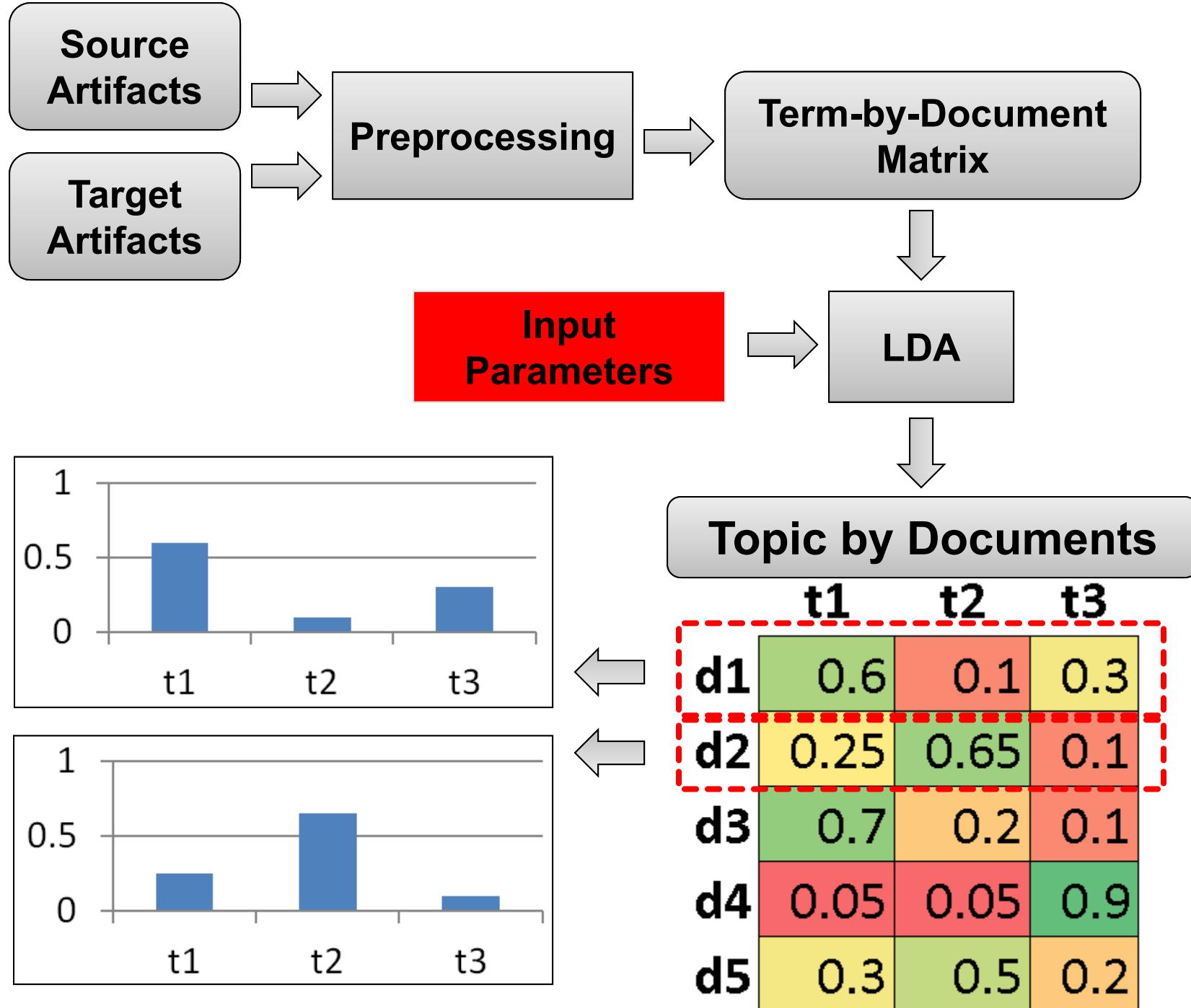


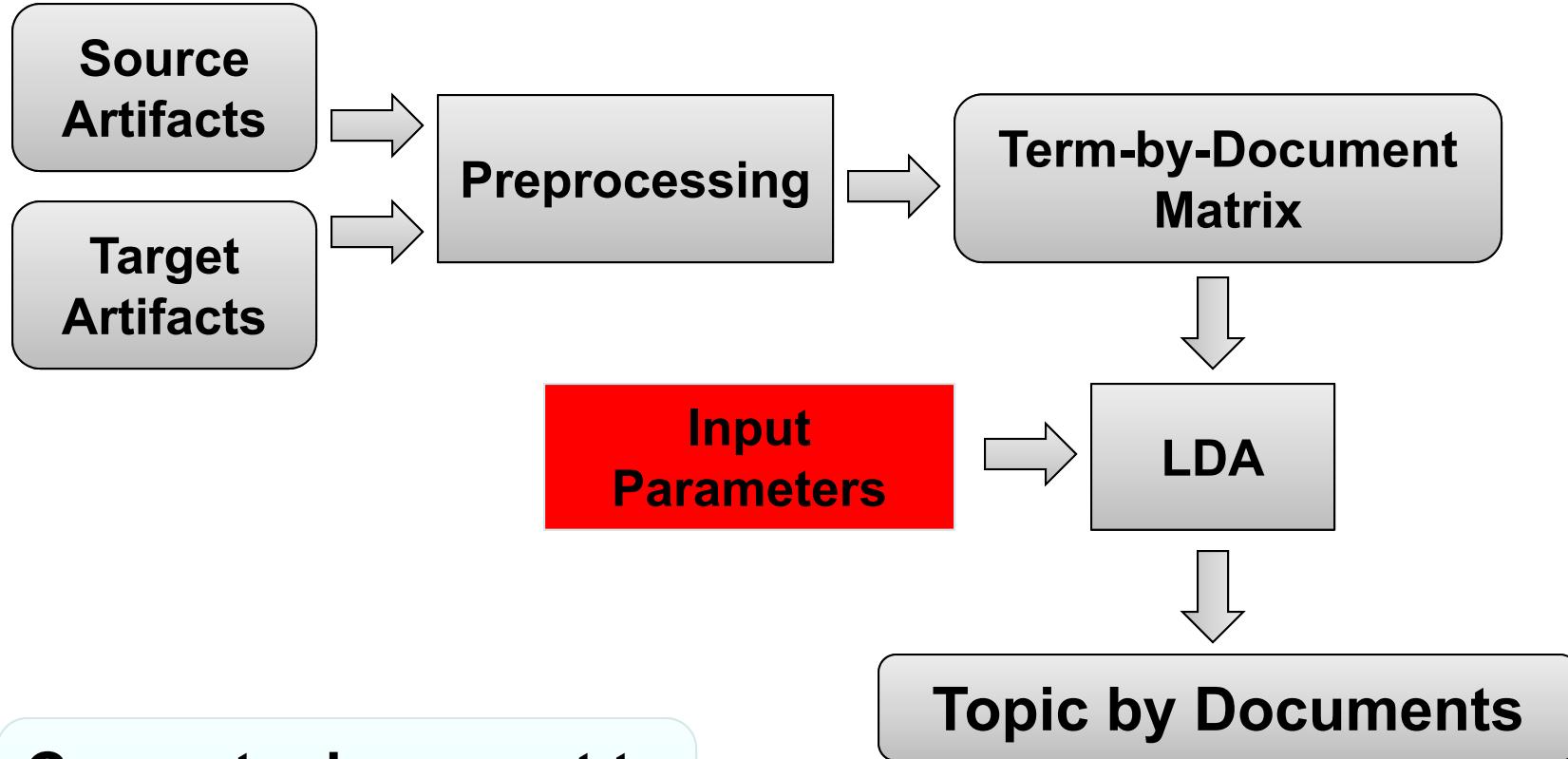
Topic by Documents

	t1	t2	t3
d1	0.6	0.1	0.3
d2	0.25	0.65	0.1
d3	0.7	0.2	0.1
d4	0.05	0.05	0.9
d5	0.3	0.5	0.2









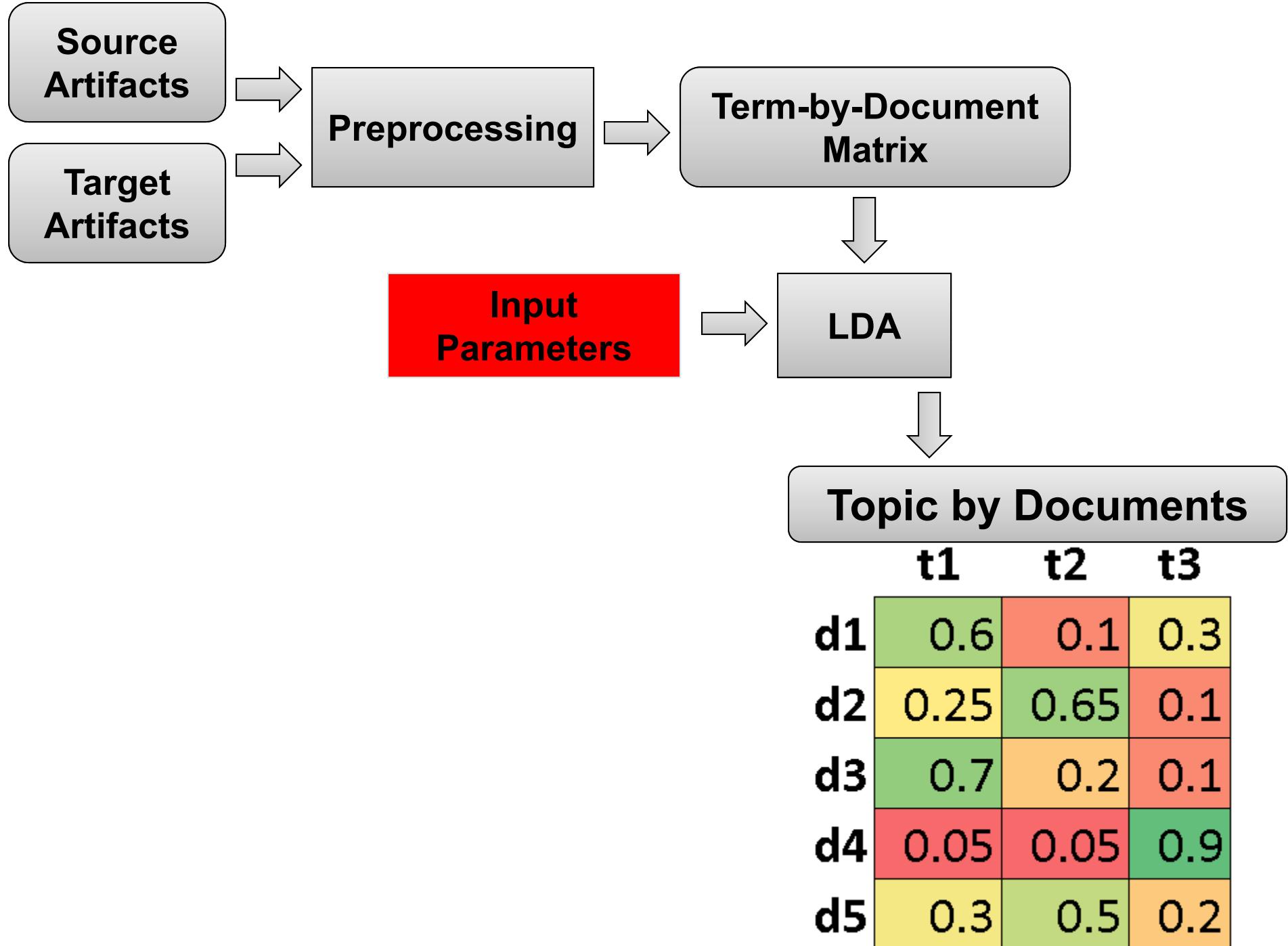
Compute document to document similarity

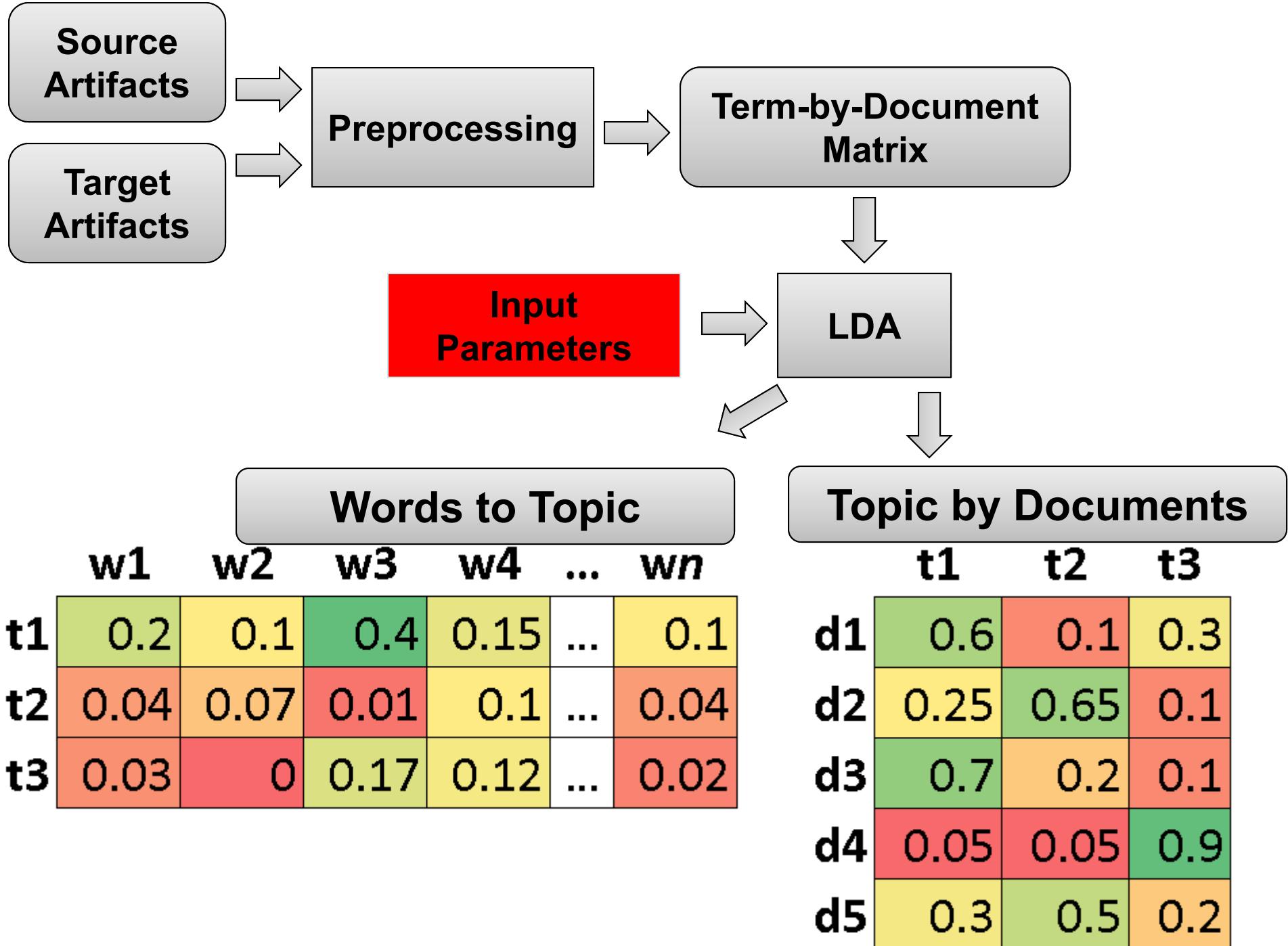
Compute query to document similarity

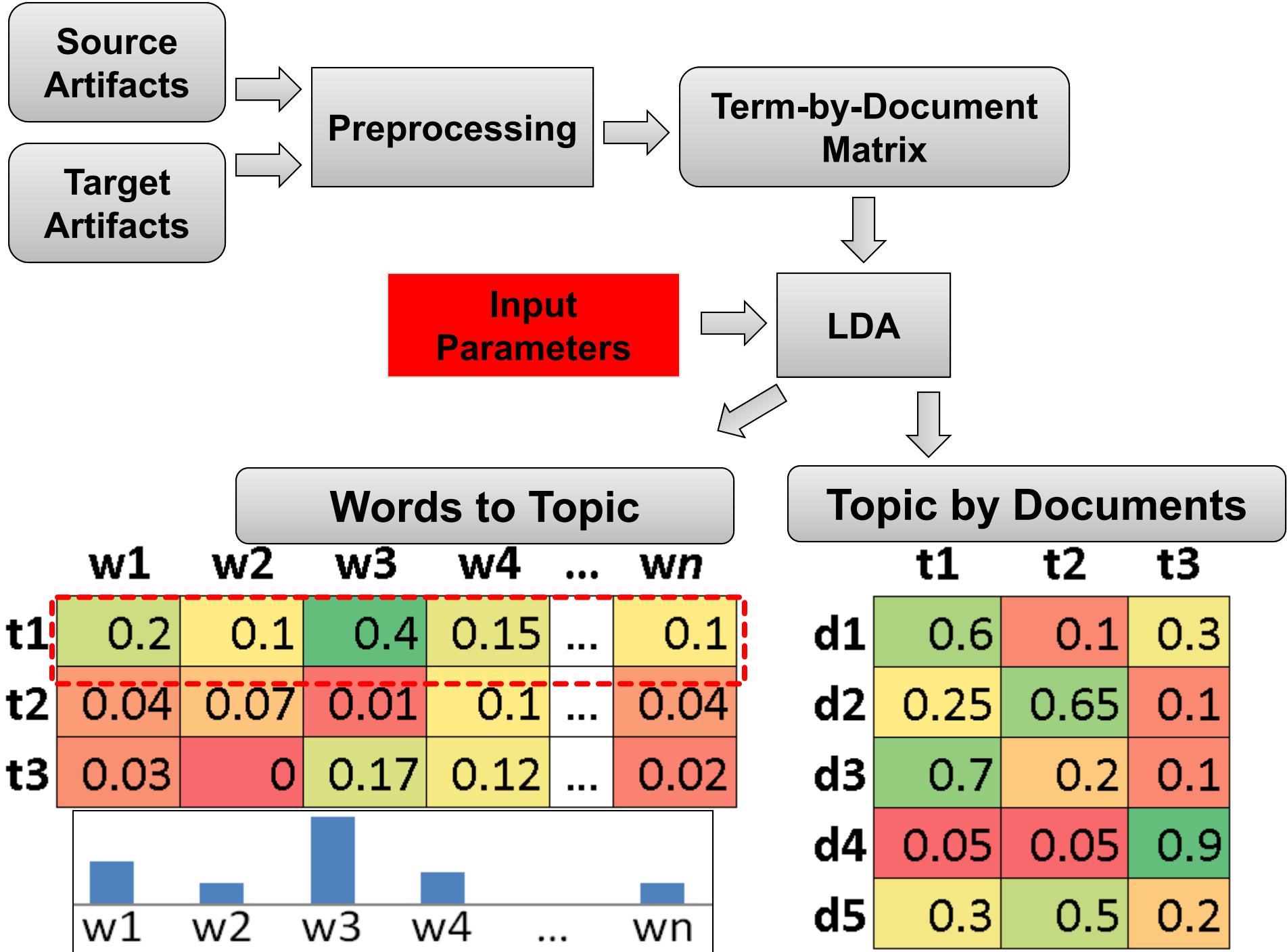
“Cluster” documents by topics

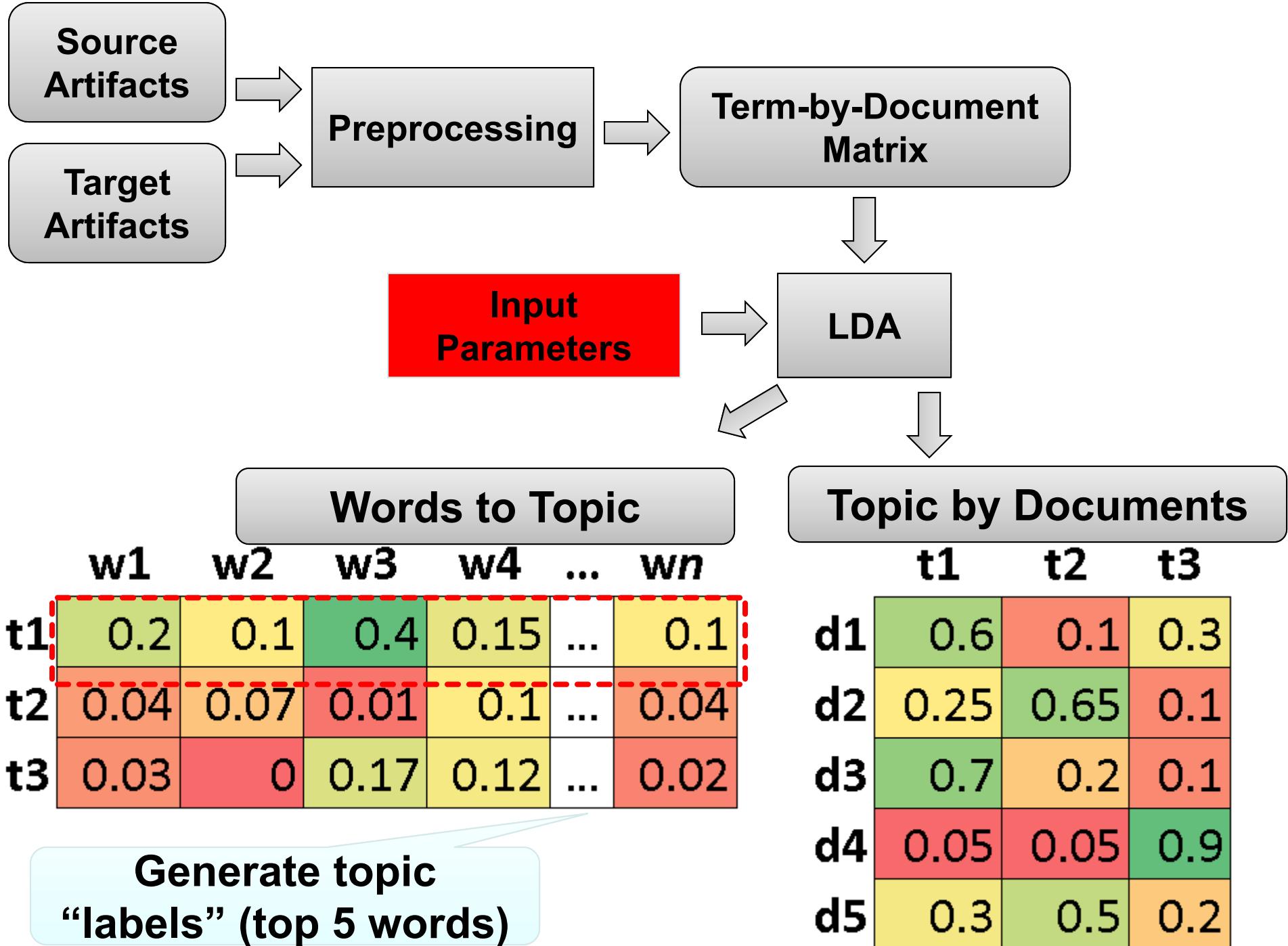
Topic by Documents

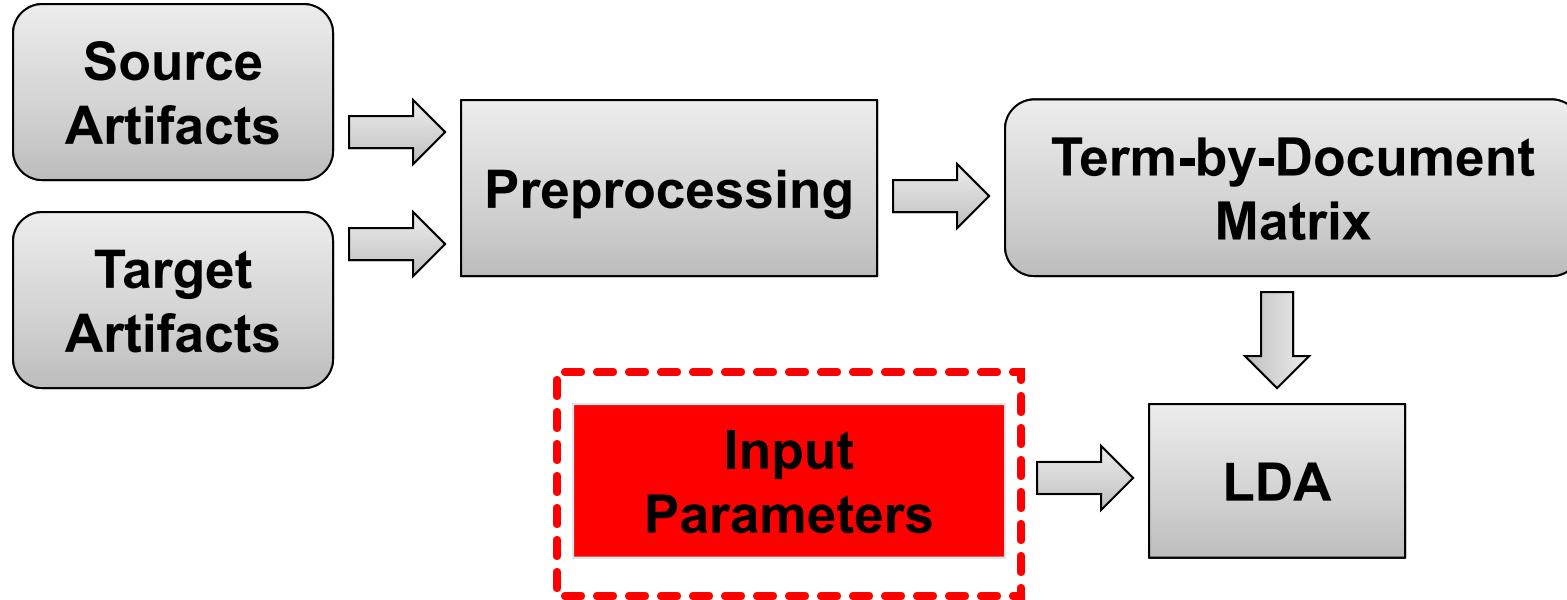
	t1	t2	t3
d1	0.6	0.1	0.3
d2	0.25	0.65	0.1
d3	0.7	0.2	0.1
d4	0.05	0.05	0.9
d5	0.3	0.5	0.2

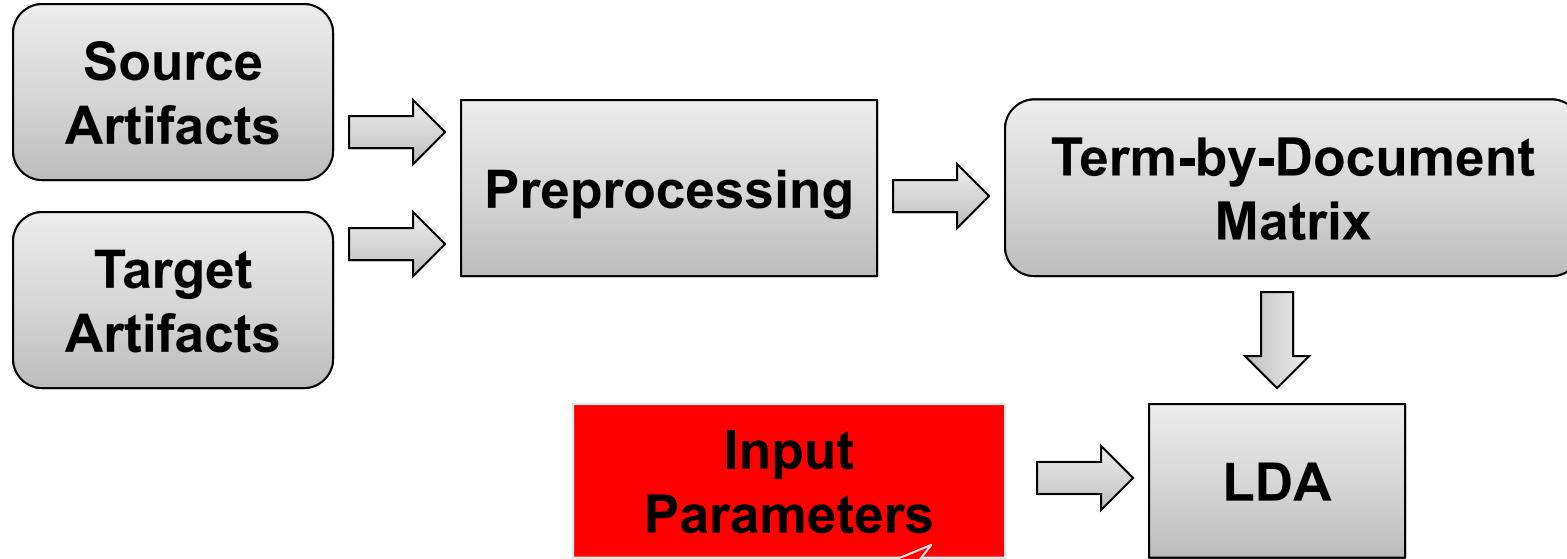












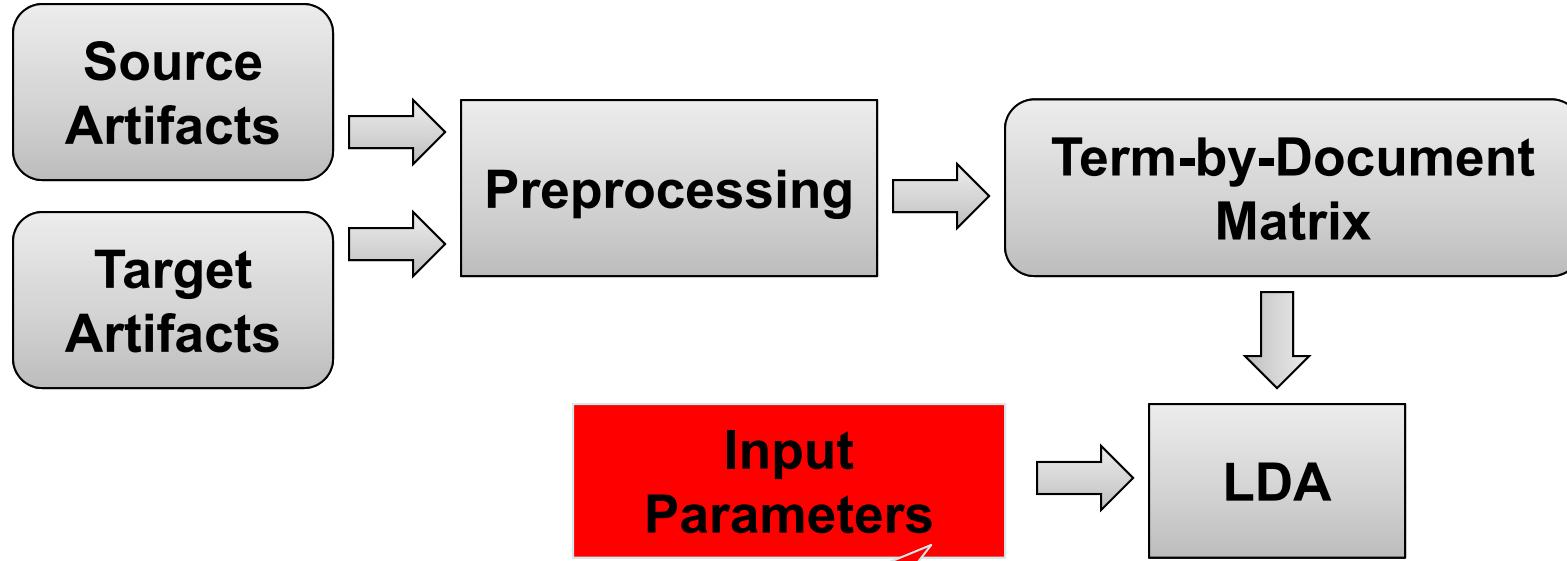
Configuration

of iterations

of topics

α

β



Configuration

of iterations

of topics

α

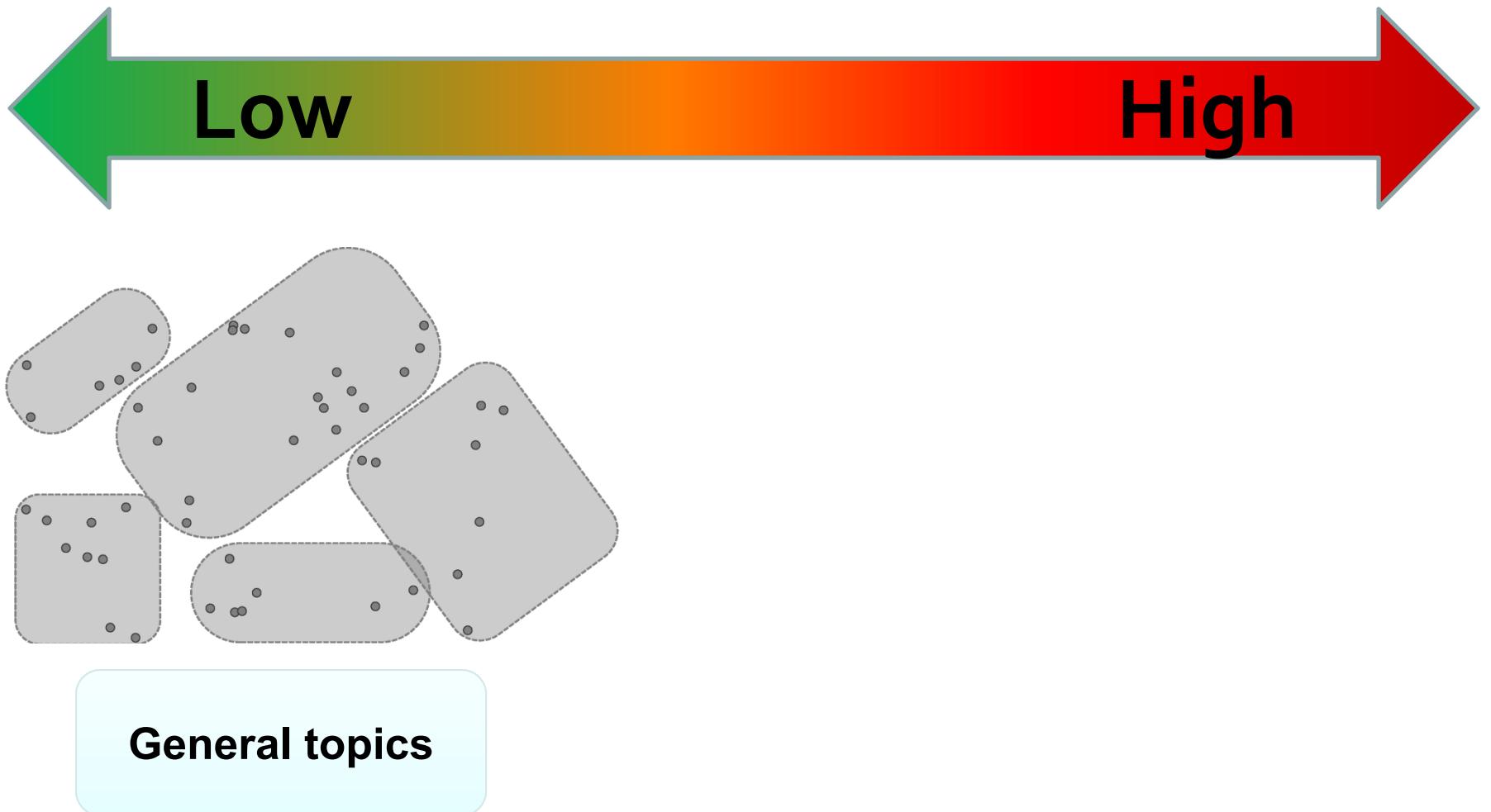
β

Number of Gibbs
samplings of LDA
model

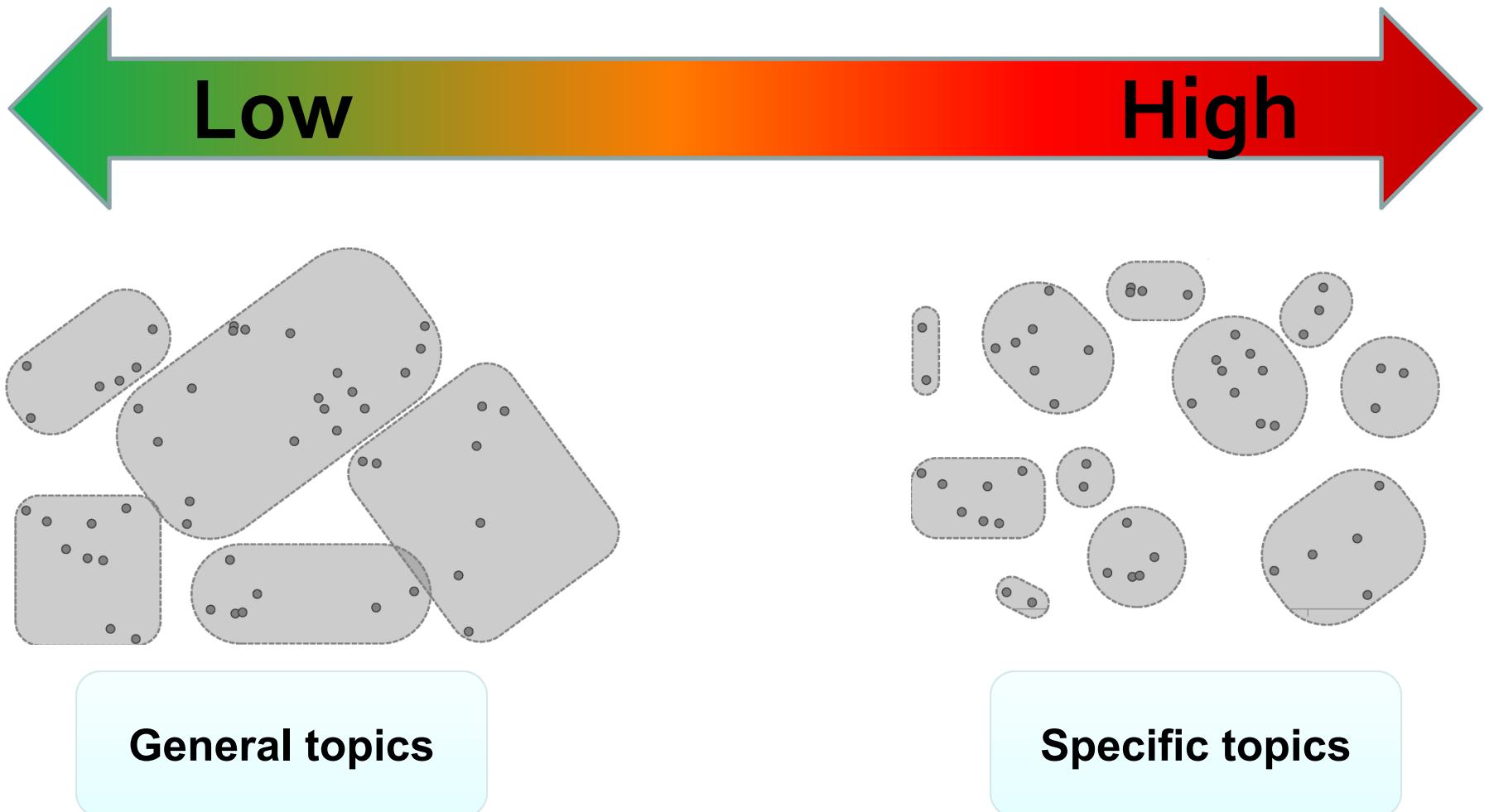
Number of topics...



Number of topics...



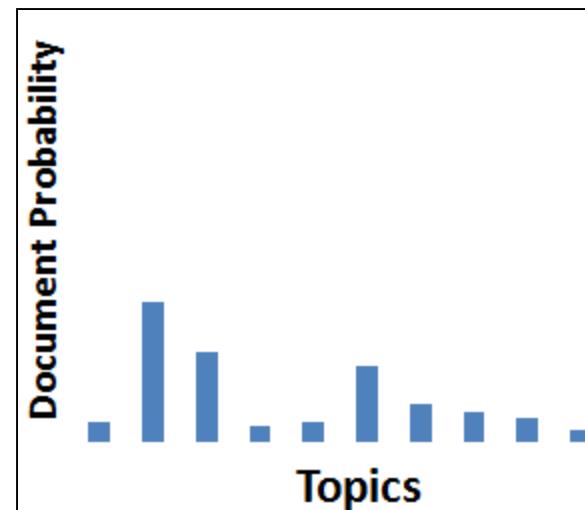
Number of topics...



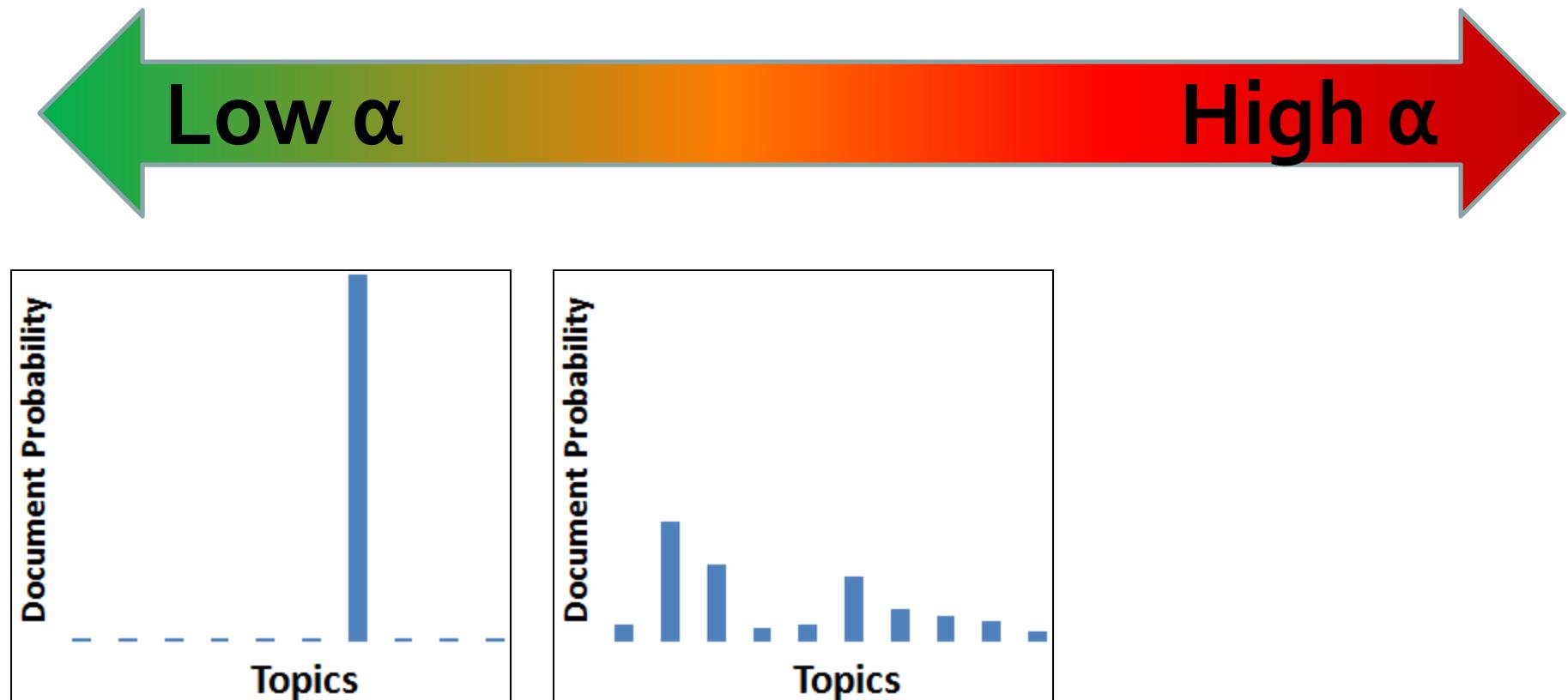
α , influences the “smoothness” of documents to topics distribution



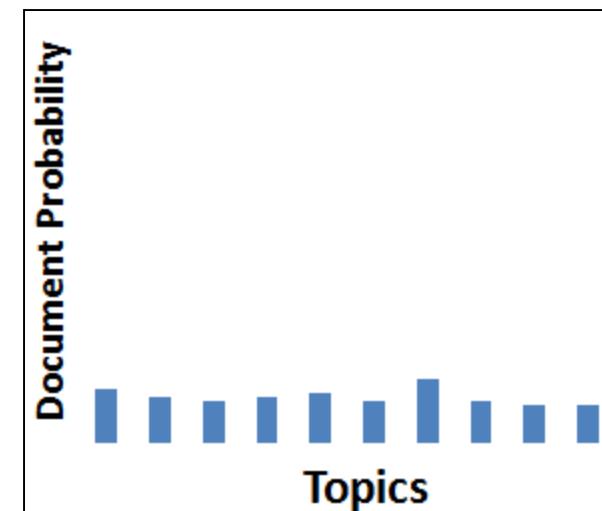
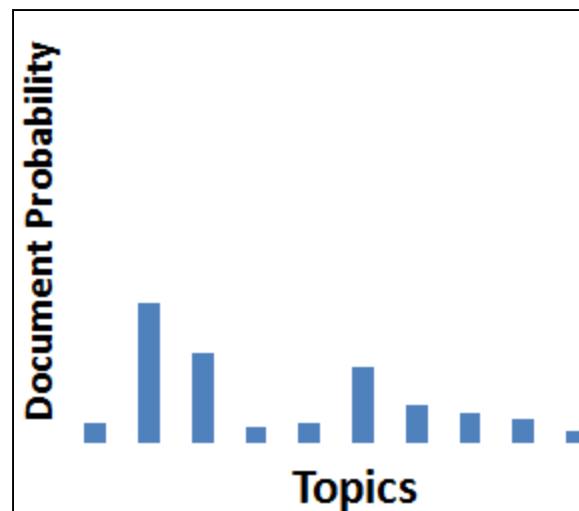
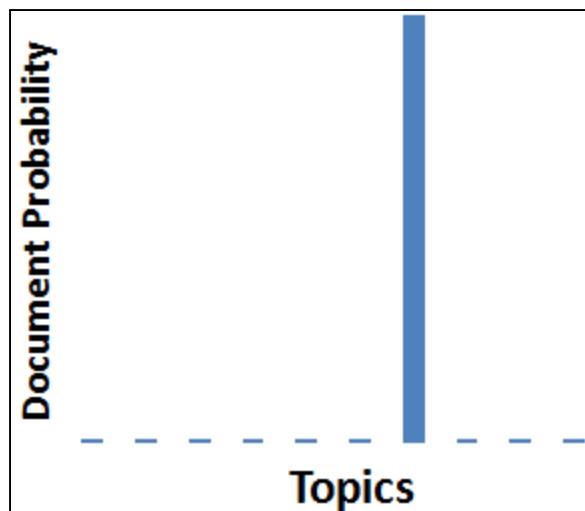
α , influences the “smoothness” of documents to topics distribution



α , influences the “smoothness” of documents to topics distribution



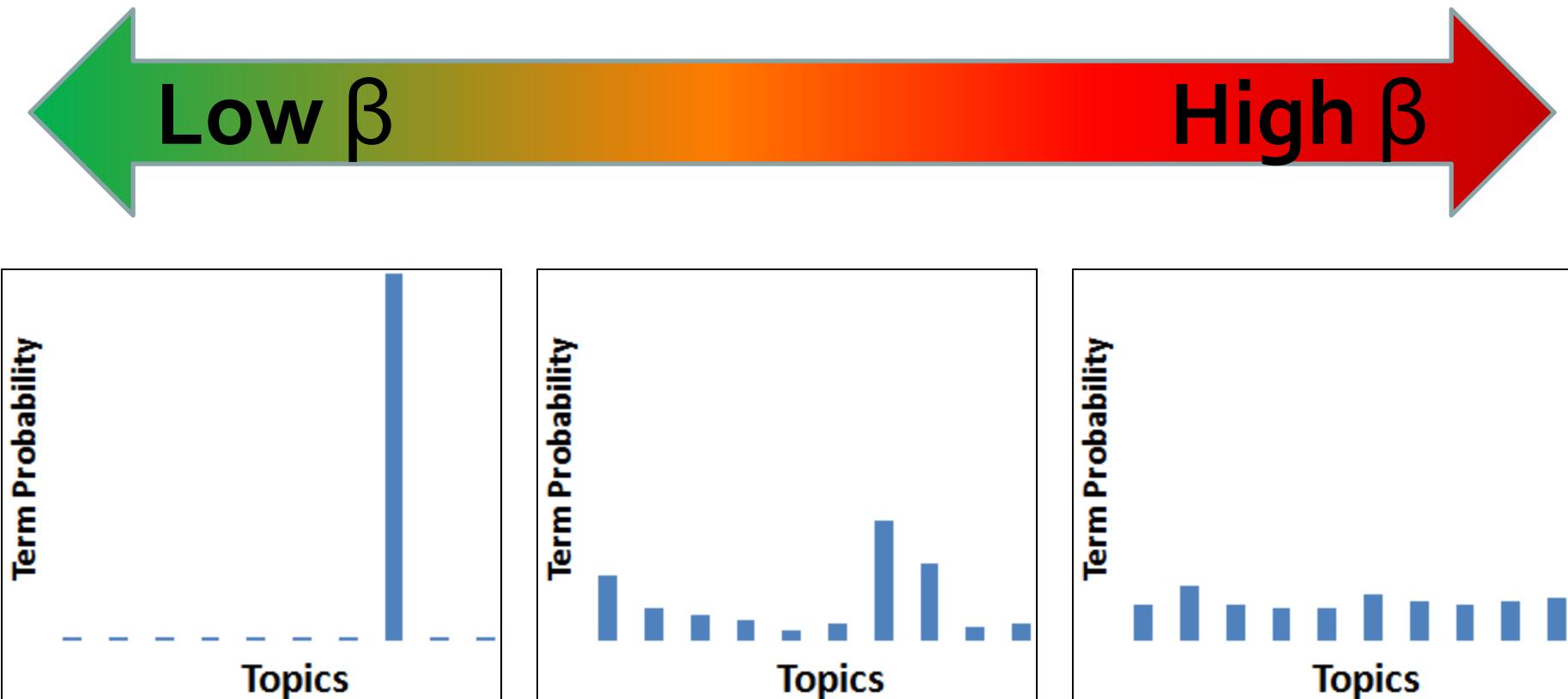
α , influences the “smoothness” of documents to topics distribution



β , influences the “smoothness” of topics to words distribution



β , influences the “smoothness” of topics to words distribution

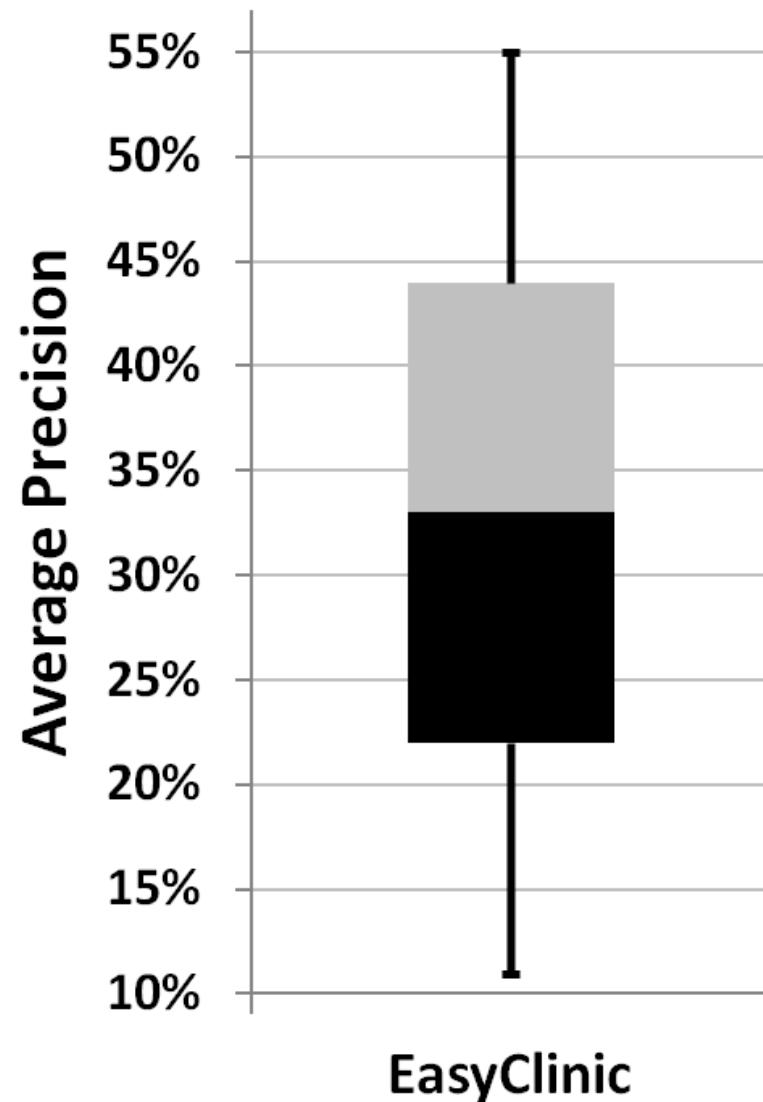


LDA parameters significantly influence the results

- 1,000 different configurations of LDA parameters
 - Evaluate the Average Precision on EasyClinic

LDA parameters significantly influence the results

- 1,000 different configurations of LDA parameters
 - Evaluate the Average Precision on EasyClinic

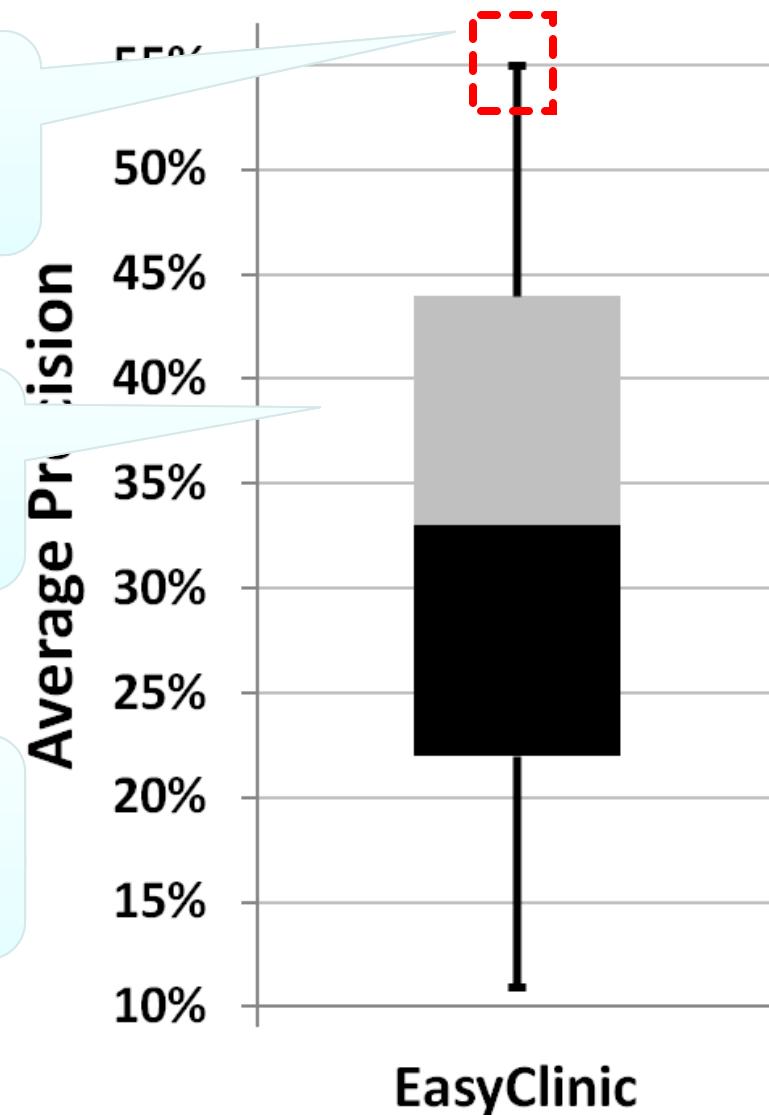


LDA parameters significantly influence the results

- Few configurations produce good results
- LDA parameters

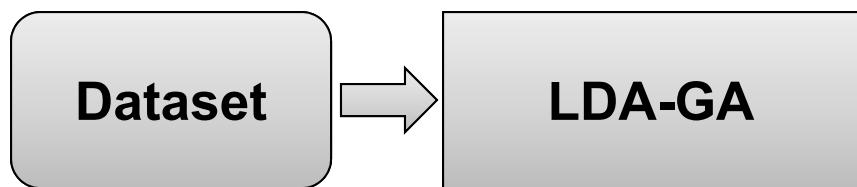
High variability in results

LDA-GA

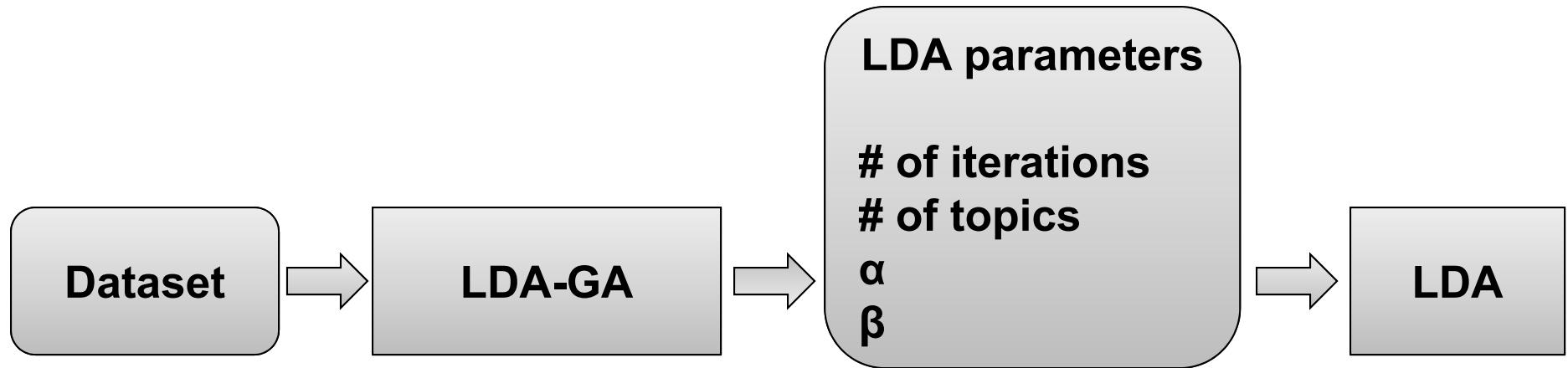


LDA-GA: automatically calibrate the input parameters of LDA using genetic algorithms

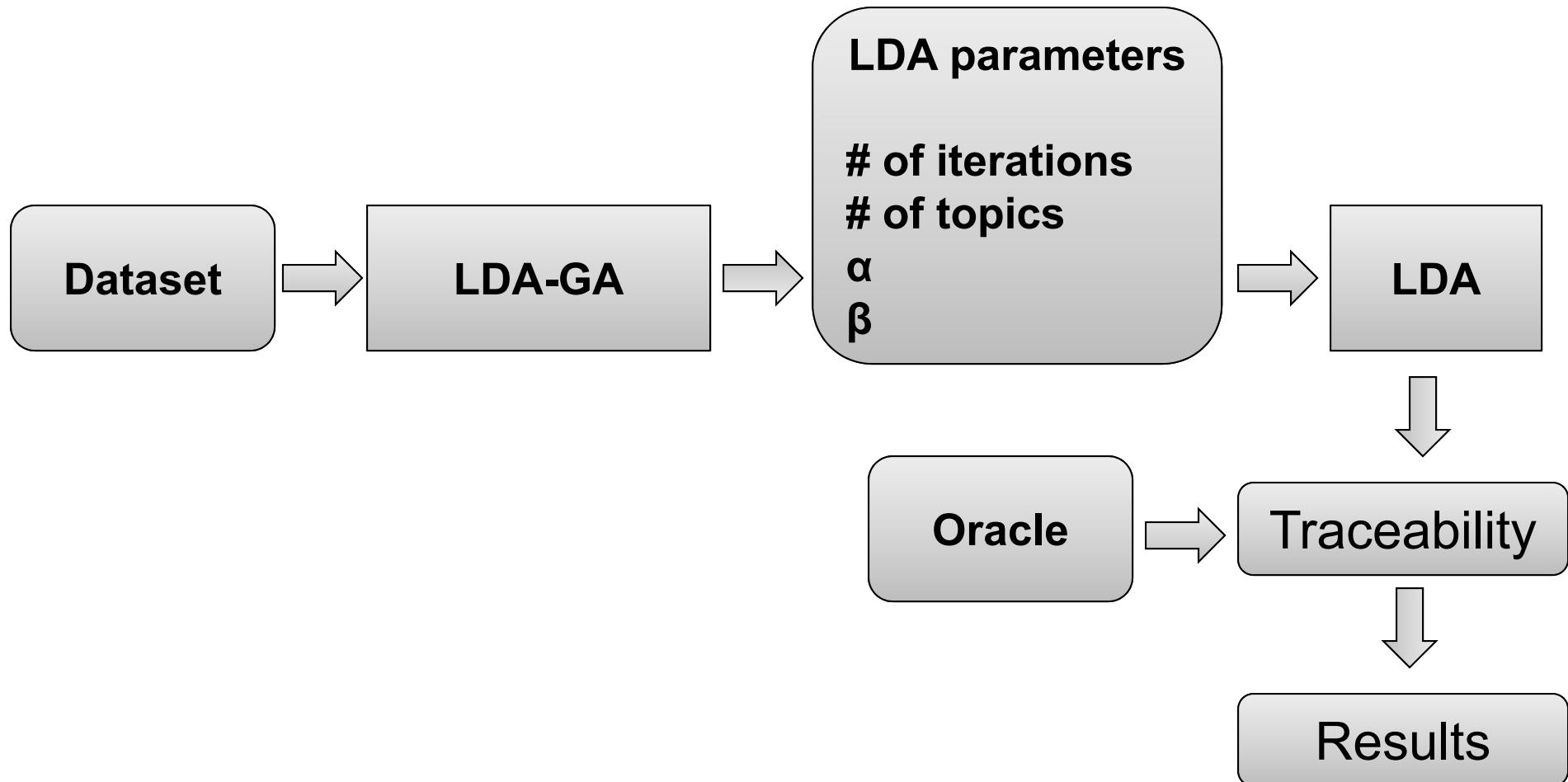
LDA-GA: automatically calibrate the input parameters of LDA using genetic algorithms



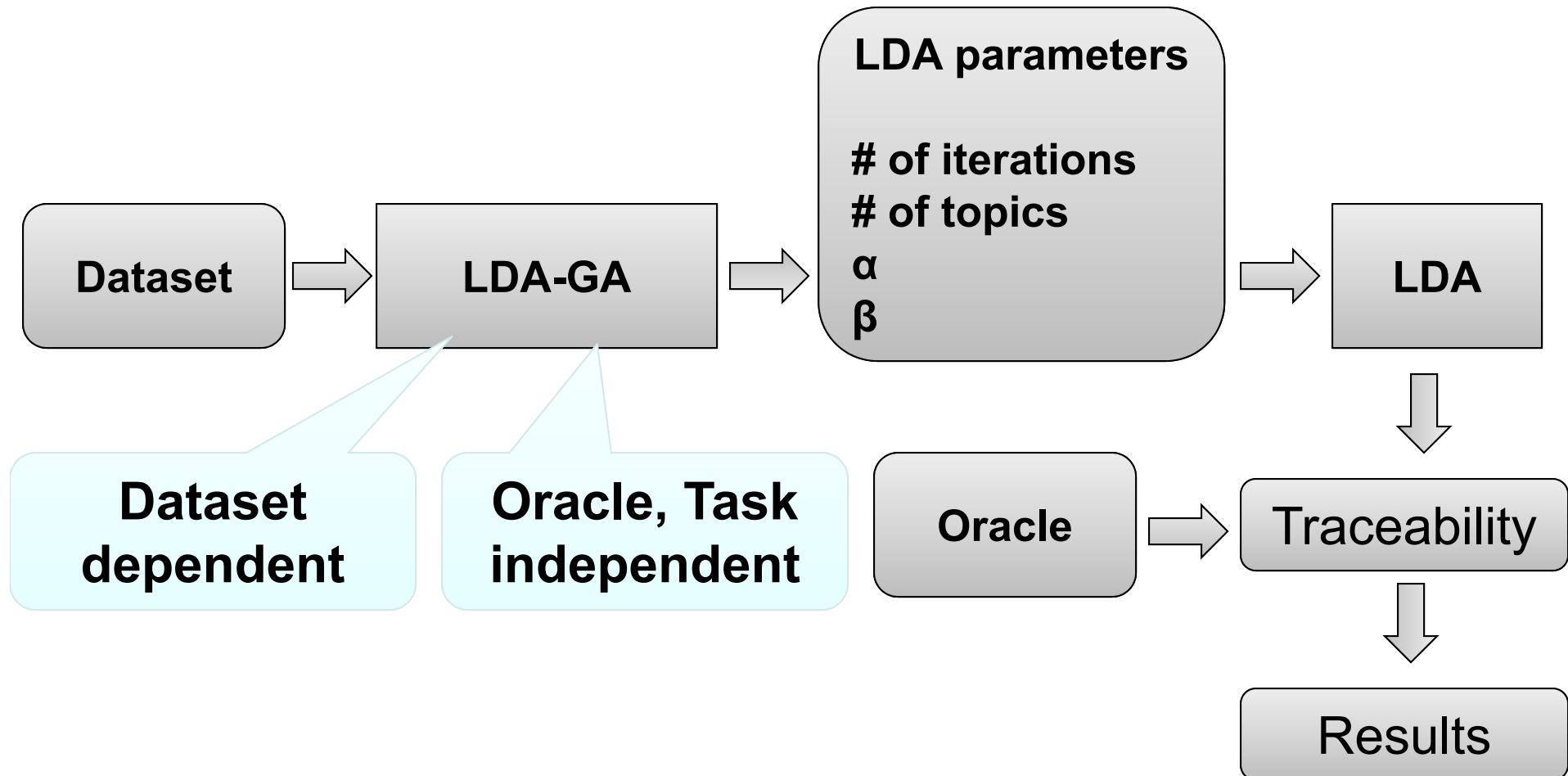
LDA-GA: automatically calibrate the input parameters of LDA using genetic algorithms



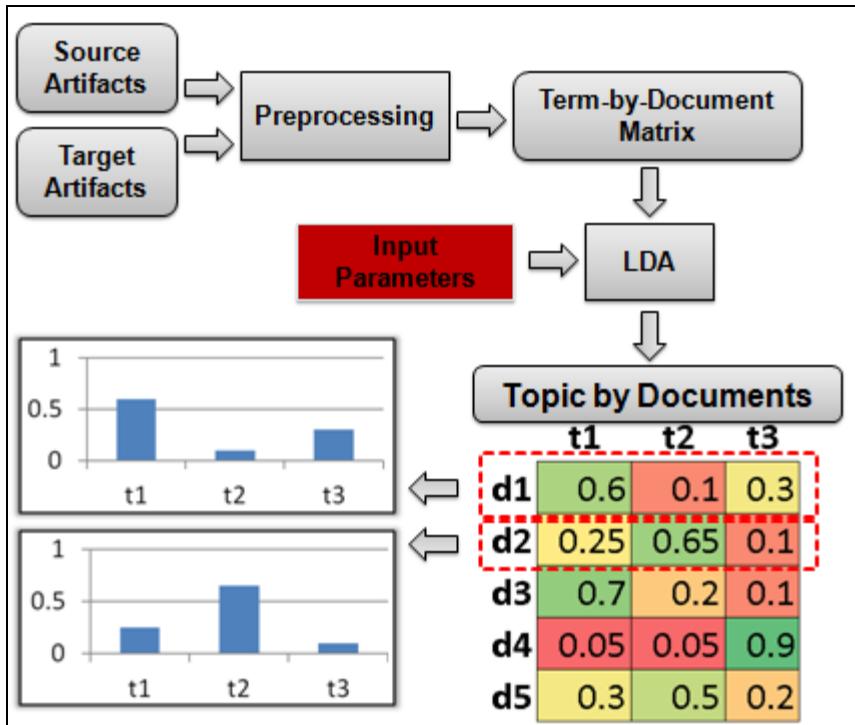
LDA-GA: automatically calibrate the input parameters of LDA using genetic algorithms



LDA-GA: automatically calibrate the input parameters of LDA using genetic algorithms



Evaluating the LDA parameter configuration



Topic by Documents

	t1	t2	t3
d1	0.6	0.1	0.3
d2	0.25	0.65	0.1
d3	0.7	0.2	0.1
d4	0.05	0.05	0.9
d5	0.3	0.5	0.2

Topic by Documents

	t1	t2	t3
d1	0.6	0.1	0.3
d2	0.25	0.65	0.1
d3	0.7	0.2	0.1
d4	0.05	0.05	0.9
d5	0.3	0.5	0.2

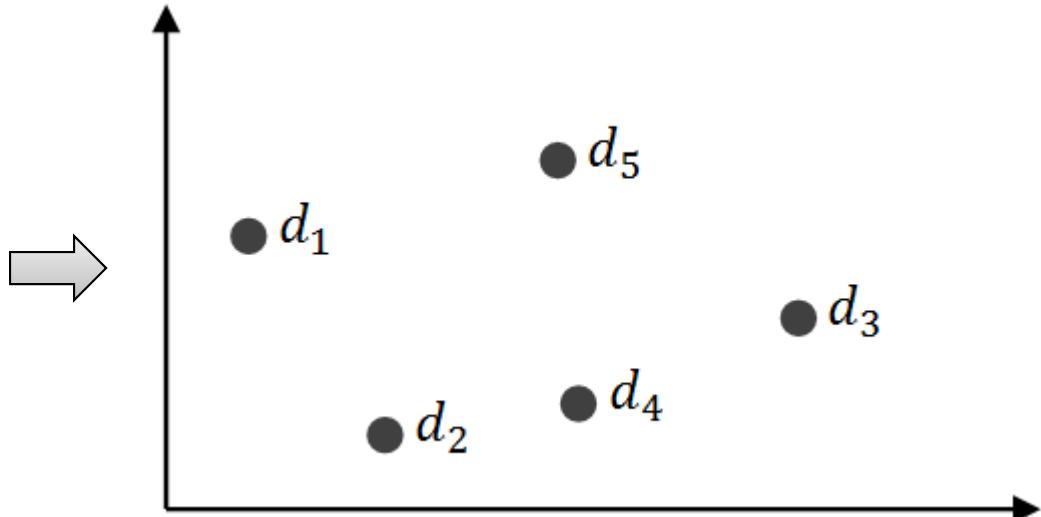
Dominant
Topics

Topic by Documents

	t1	t2	t3
d1	0.6	0.1	0.3
d2	0.25	0.65	0.1
d3	0.7	0.2	0.1
d4	0.05	0.05	0.9
d5	0.3	0.5	0.2

Dominant
Topics

LDA Model

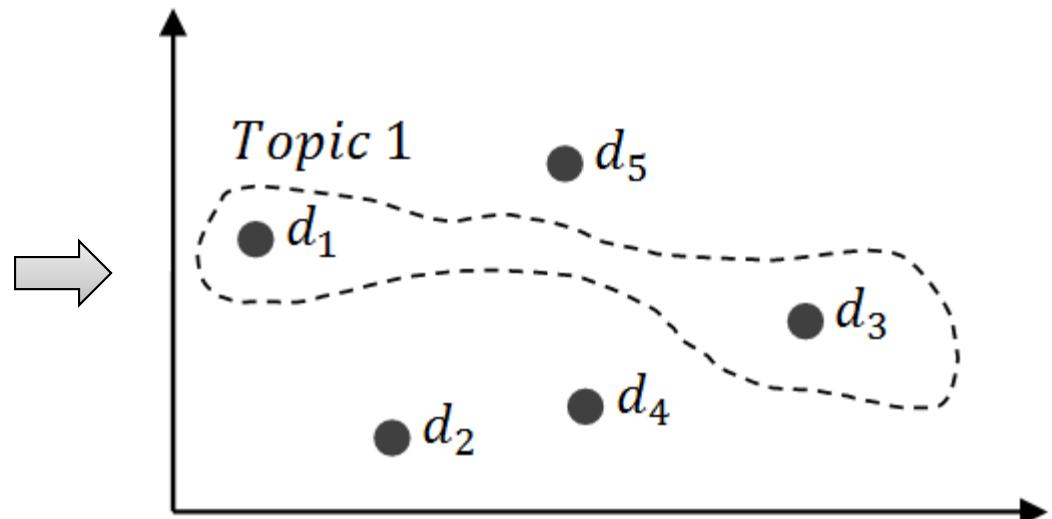


Topic by Documents

	t1	t2	t3
d1	0.6	0.1	0.3
d2	0.25	0.65	0.1
d3	0.7	0.2	0.1
d4	0.05	0.05	0.9
d5	0.3	0.5	0.2

Dominant
Topics

LDA Model

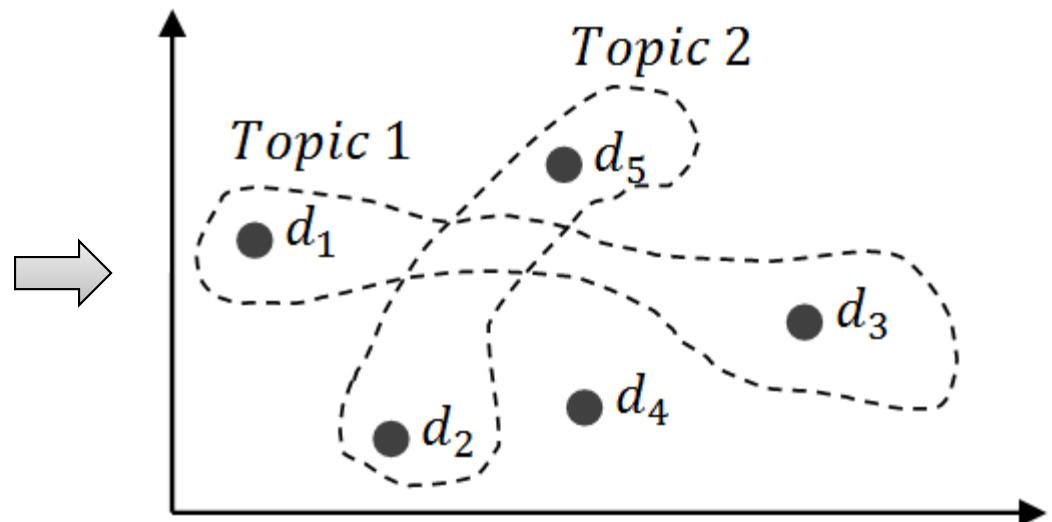


Topic by Documents

	t1	t2	t3
d1	0.6	0.1	0.3
d2	0.25	0.65	0.1
d3	0.7	0.2	0.1
d4	0.05	0.05	0.9
d5	0.3	0.5	0.2

Dominant
Topics

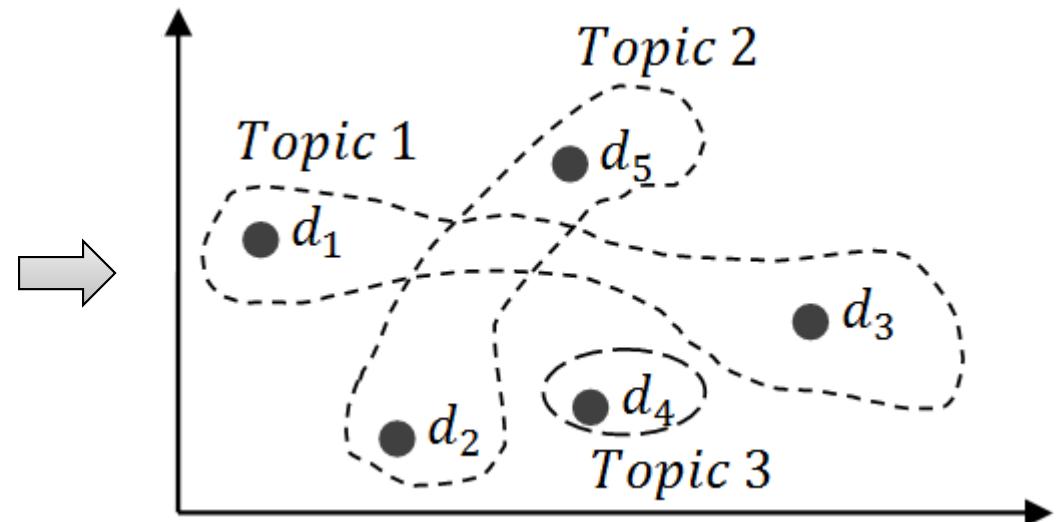
LDA Model



Topic by Documents

	t1	t2	t3
d1	0.6	0.1	0.3
d2	0.25	0.65	0.1
d3	0.7	0.2	0.1
d4	0.05	0.05	0.9
d5	0.3	0.5	0.2

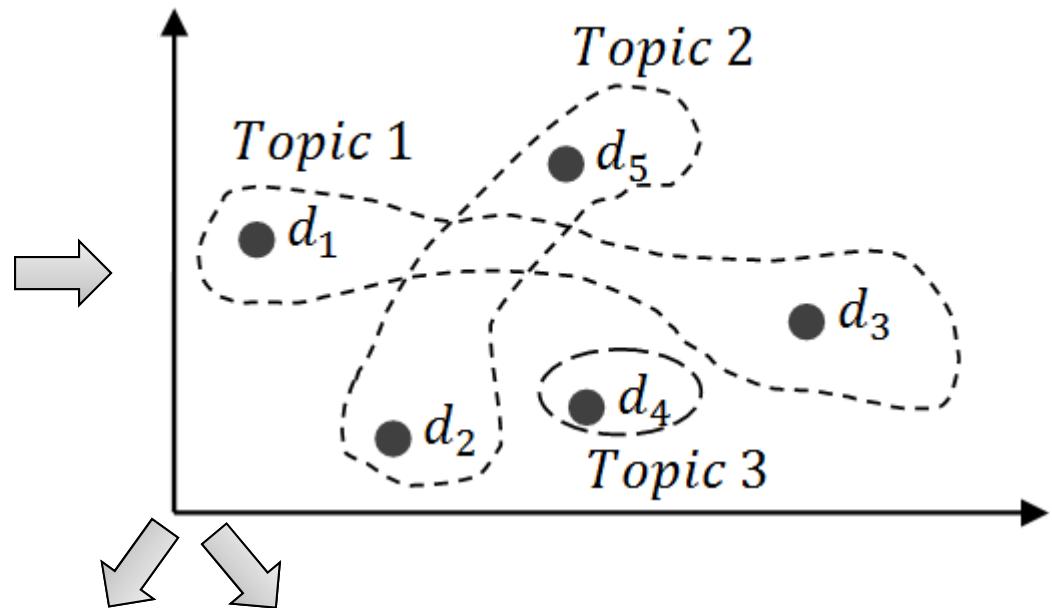
LDA Model



Topic by Documents

	t1	t2	t3
d1	0.6	0.1	0.3
d2	0.25	0.65	0.1
d3	0.7	0.2	0.1
d4	0.05	0.05	0.9
d5	0.3	0.5	0.2

LDA Model



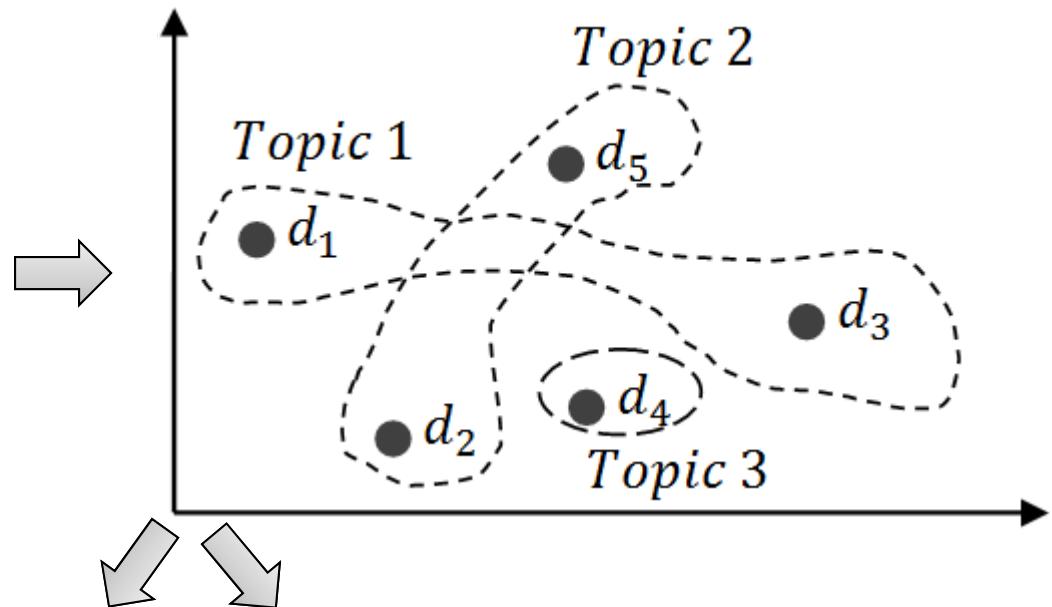
Cohesion (similarity): how related the documents in the same clusters are

Separation (dissimilarity): how distinct a cluster is from other clusters

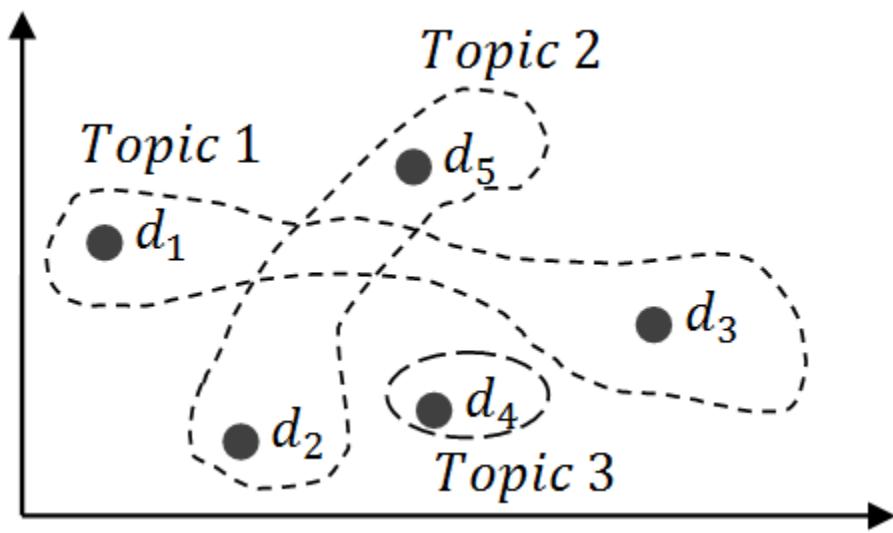
Topic by Documents

	t1	t2	t3
d1	0.6	0.1	0.3
d2	0.25	0.65	0.1
d3	0.7	0.2	0.1
d4	0.05	0.05	0.9
d5	0.3	0.5	0.2

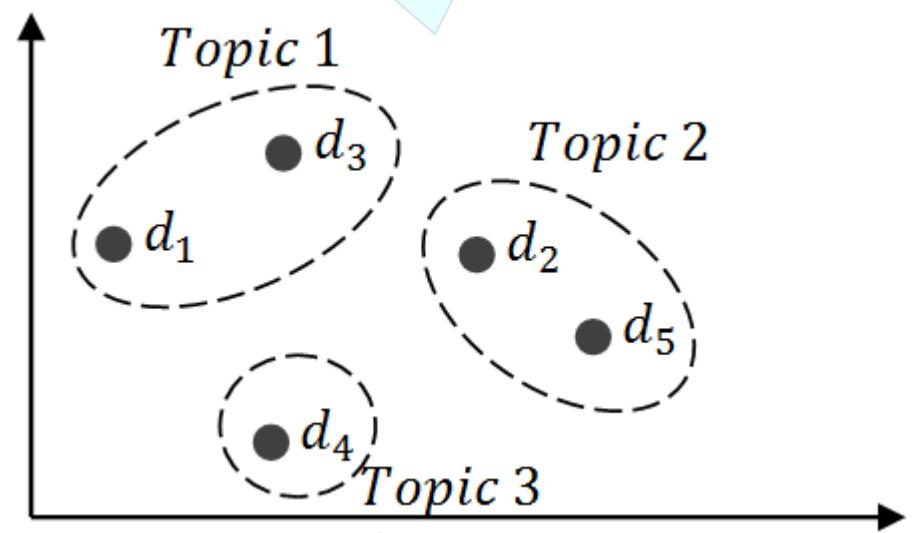
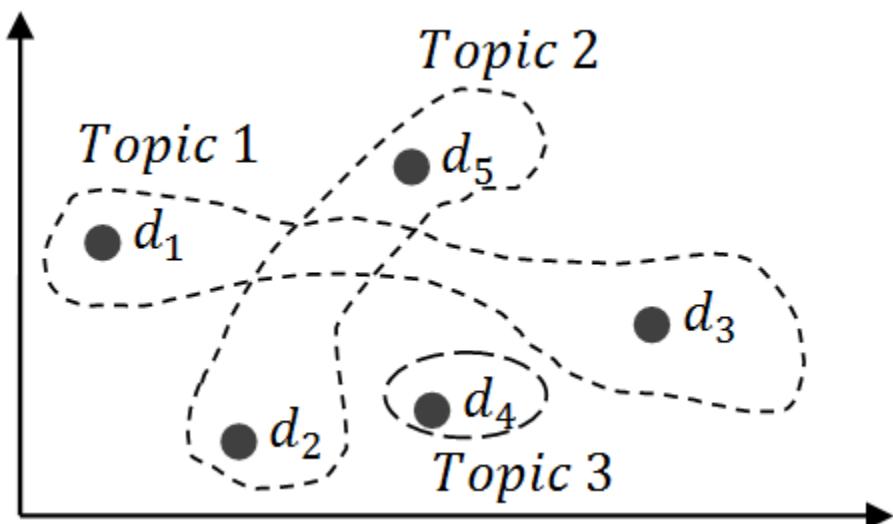
LDA Model



Silhouette coefficient



**Cohesive
Separated**



**Documents are rearranged
more closely in the n-
dimensional space**

We propose a genetic algorithm to
find those solutions

**Choose a set of Random
LDA parameters**

**Choose a set of Random
LDA parameters**

	iteration	topics	α	β
LDA Cfg. 1	510	74	0.34	2.5
LDA Cfg. 2	725	128	1.28	0.4
...
LDA Cfg. n	618	250	1.14	0.74

Choose a set of Random LDA parameters



LDA

Determine cohesion and separation

	iteration	topics	α	β
LDA Cfg. 1	510	74	0.34	2.5
LDA Cfg. 2	725	128	1.28	0.4
...
LDA Cfg. n	618	250	1.14	0.74

Choose a set of Random LDA parameters

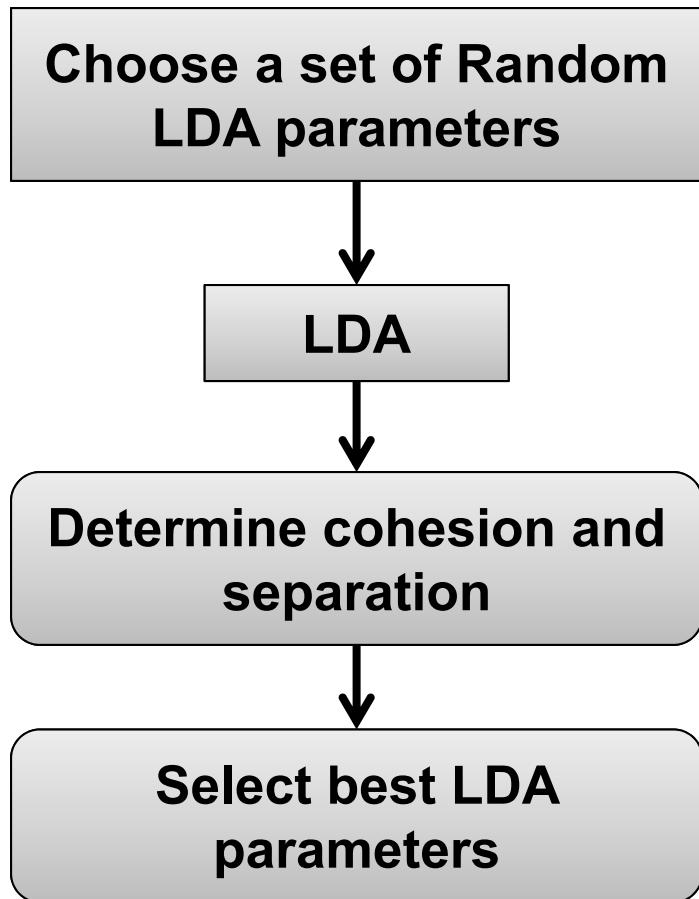


LDA

Determine cohesion and separation

	iteration	topics	α	β	Fitness
LDA Cfg. 1	510	74	0.34	2.5	0.2
LDA Cfg. 2	725	128	1.28	0.4	0.3
...	
LDA Cfg. n	618	250	1.14	0.74	0.1

Cohesion and Separation (Silhouette)



	iteration	topics	α	β	Fitness
LDA Cfg. 1	510	74	0.34	2.5	0.2
LDA Cfg. 2	725	128	1.28	0.4	0.3
...	
LDA Cfg. n	618	250	1.14	0.74	0.1

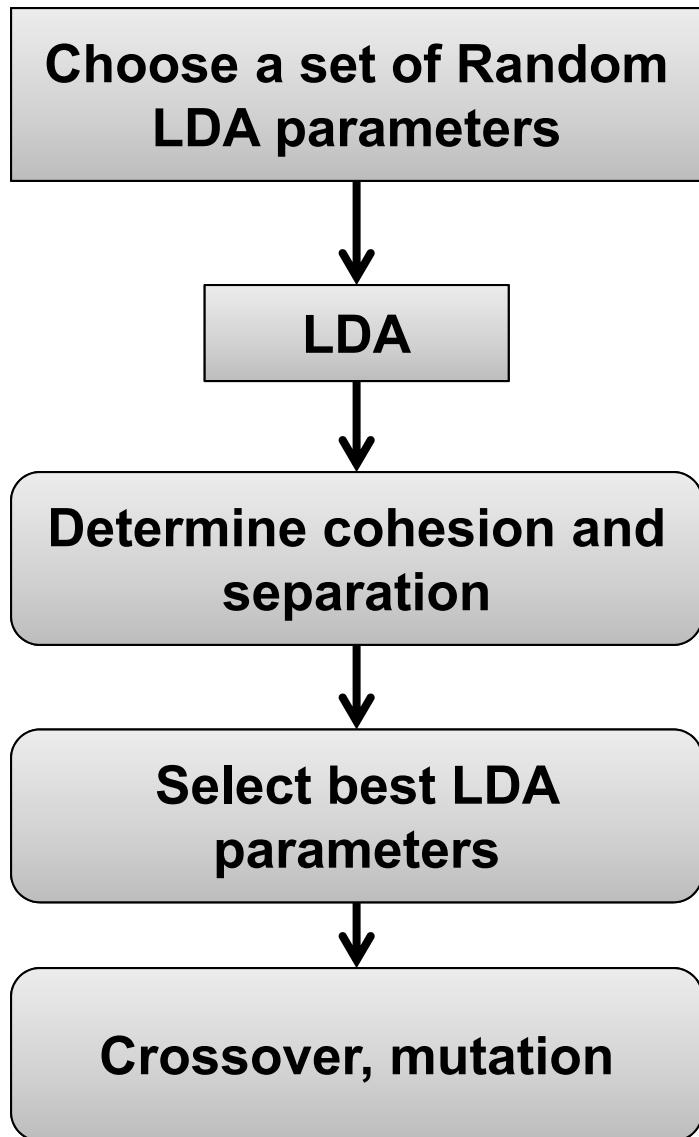
Choose a set of Random LDA parameters

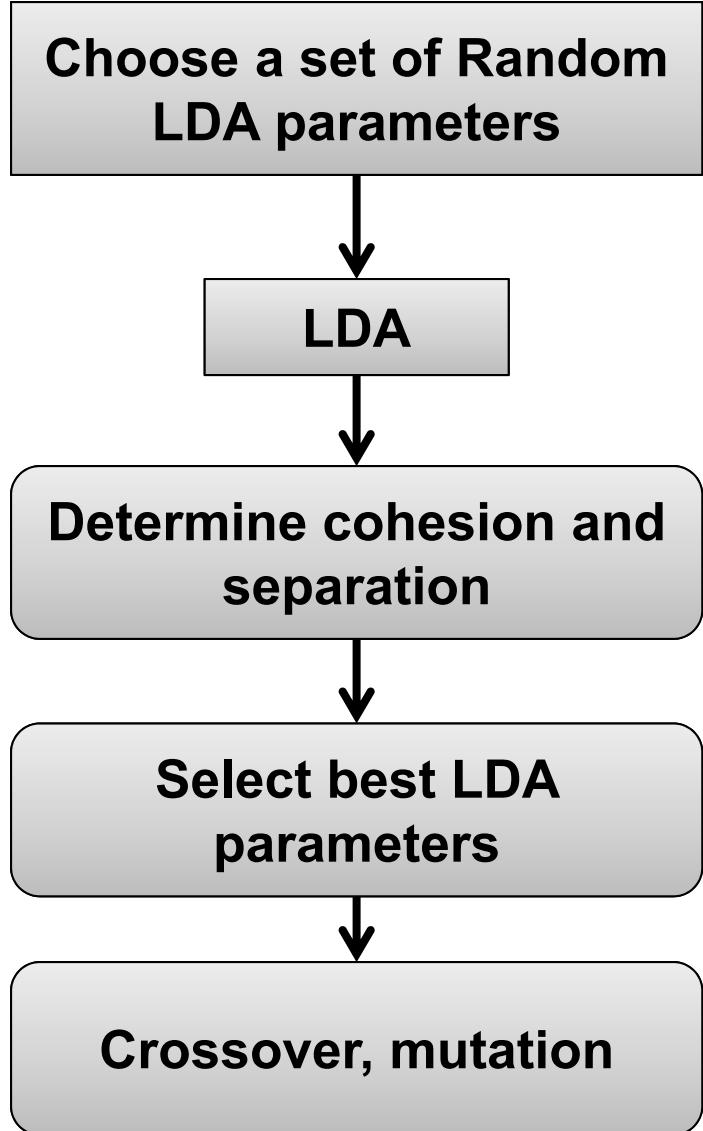
LDA

Determine cohesion and separation

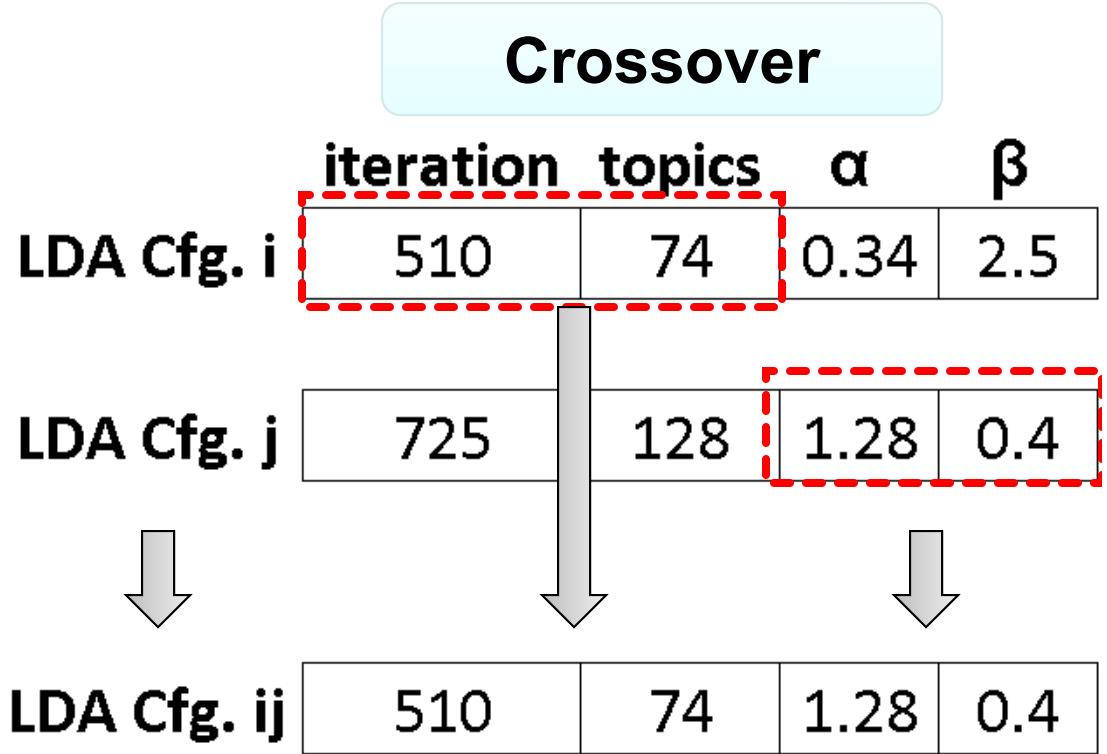
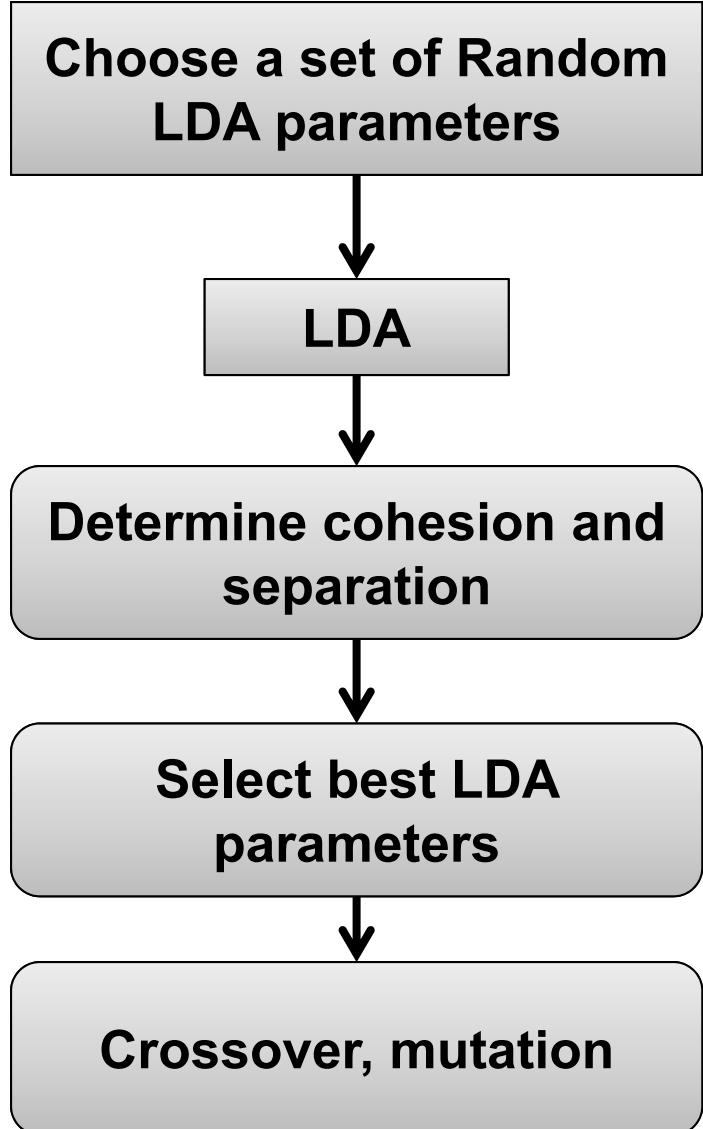
Select best LDA parameters

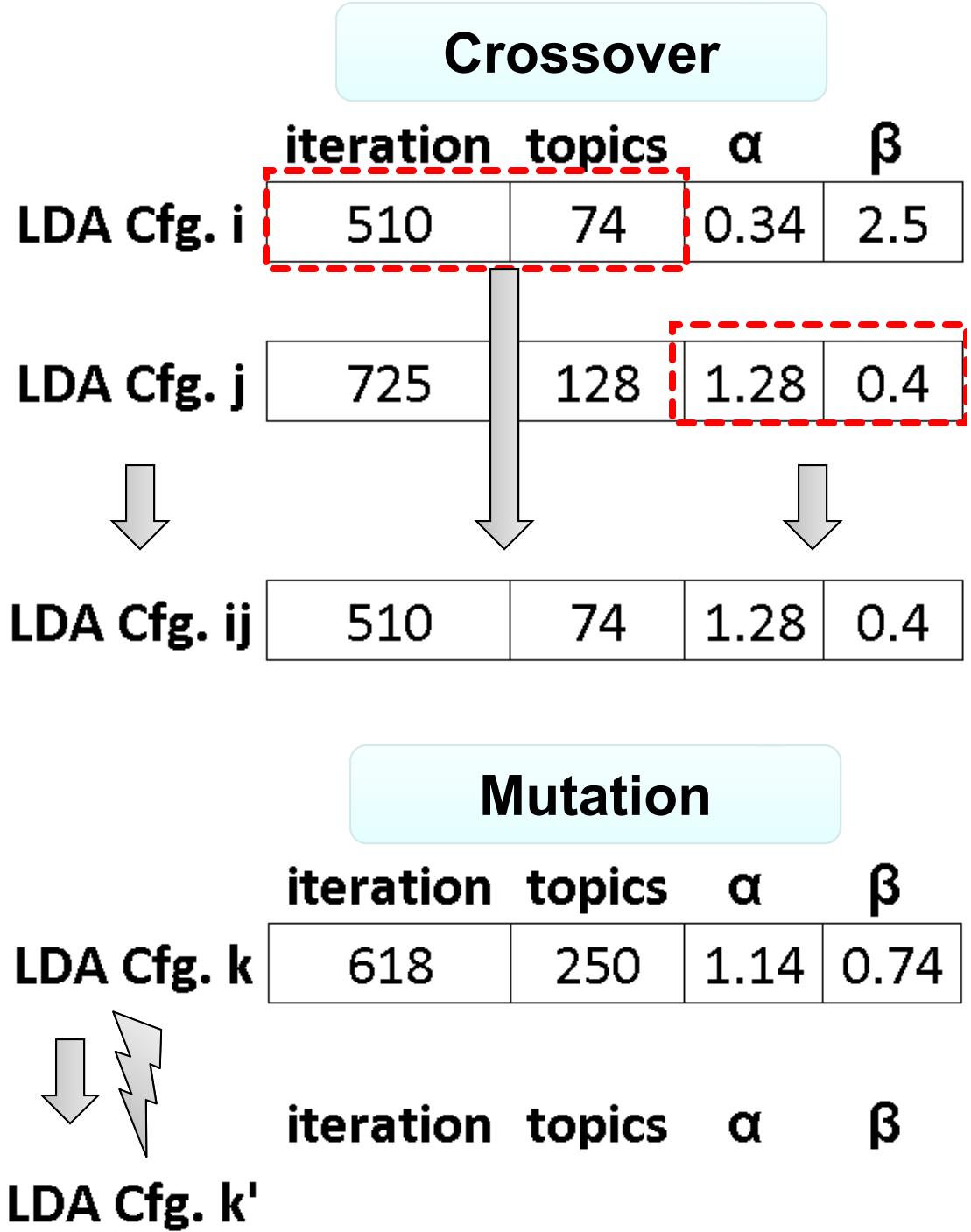
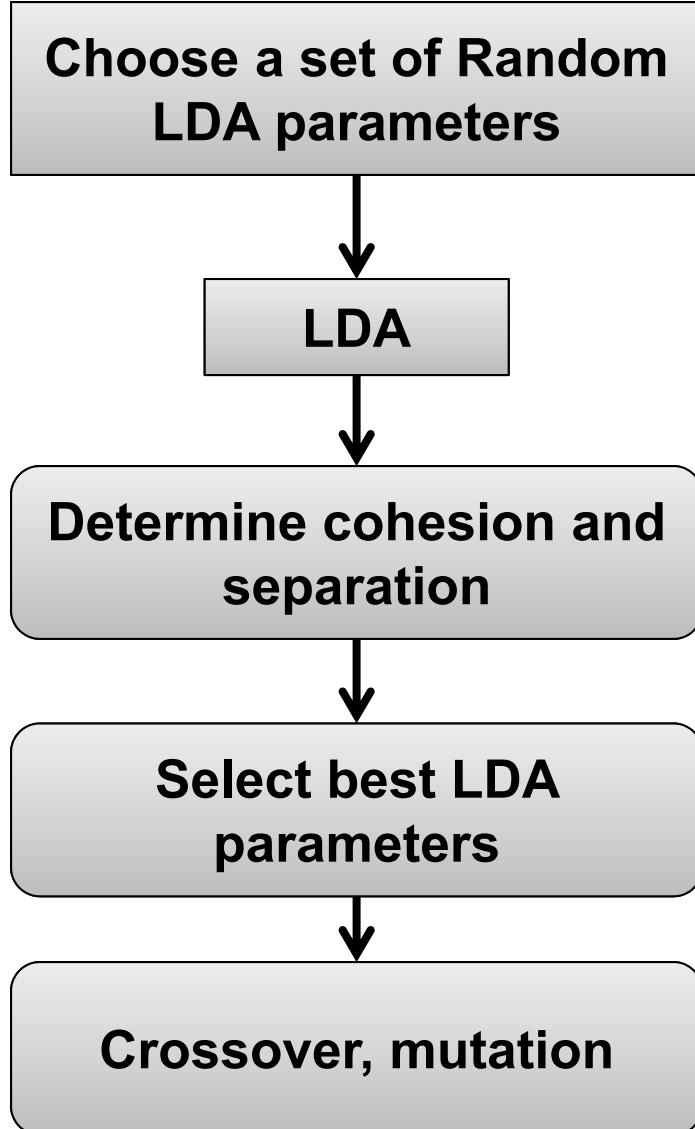
	iteration	topics	α	β	Fitness
LDA Cfg. 1	510	74	0.34	2.5	0.2
LDA Cfg. 2	725	128	1.28	0.4	0.3
...	
LDA Cfg. n	618	250	1.14	0.74	0.1

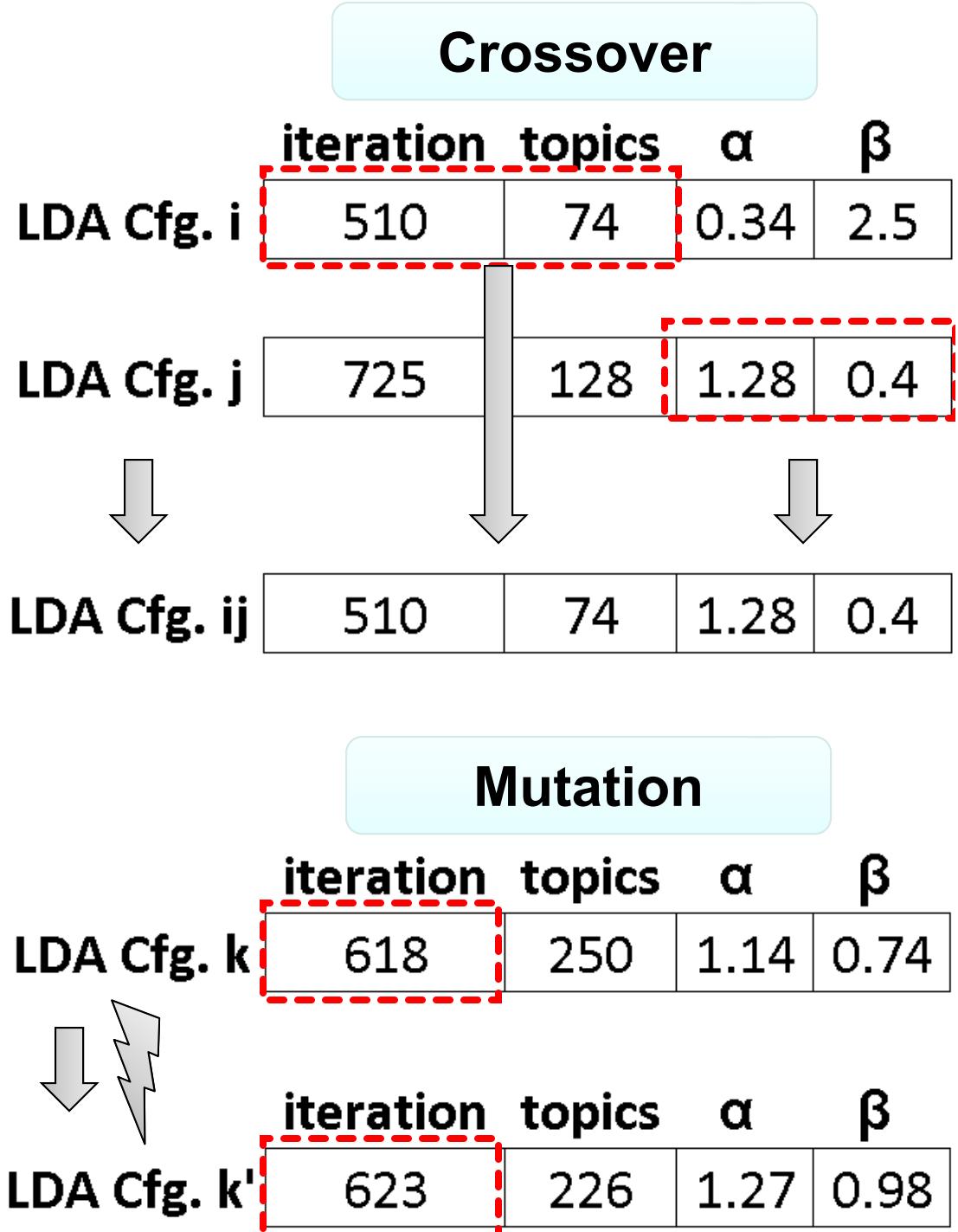
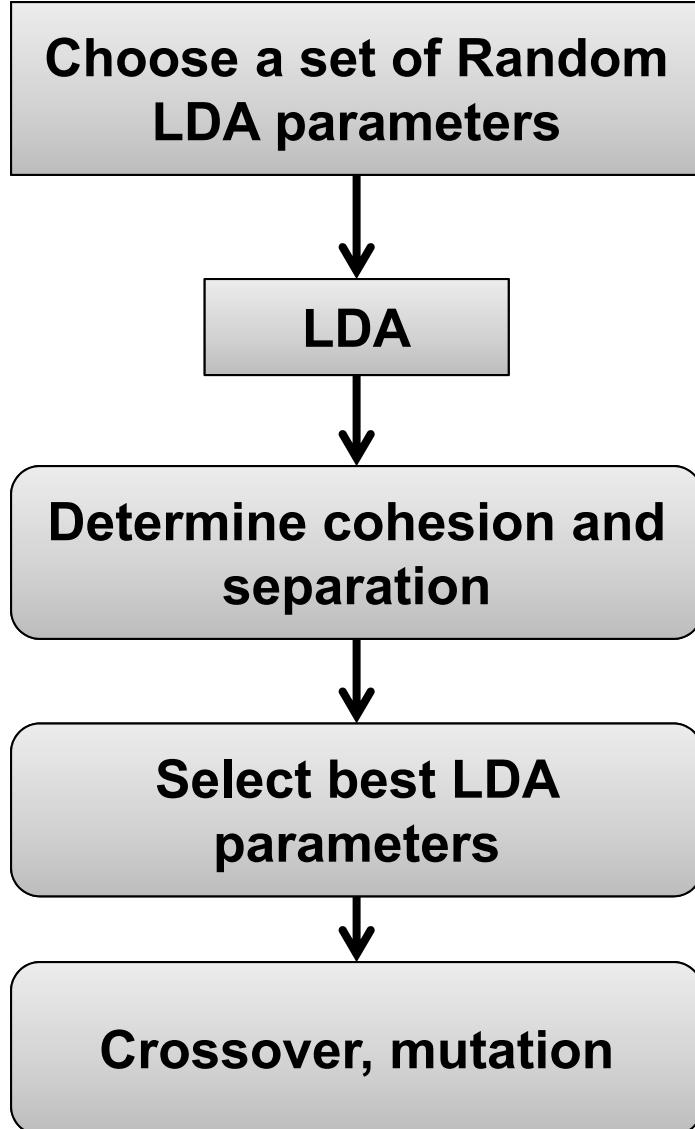


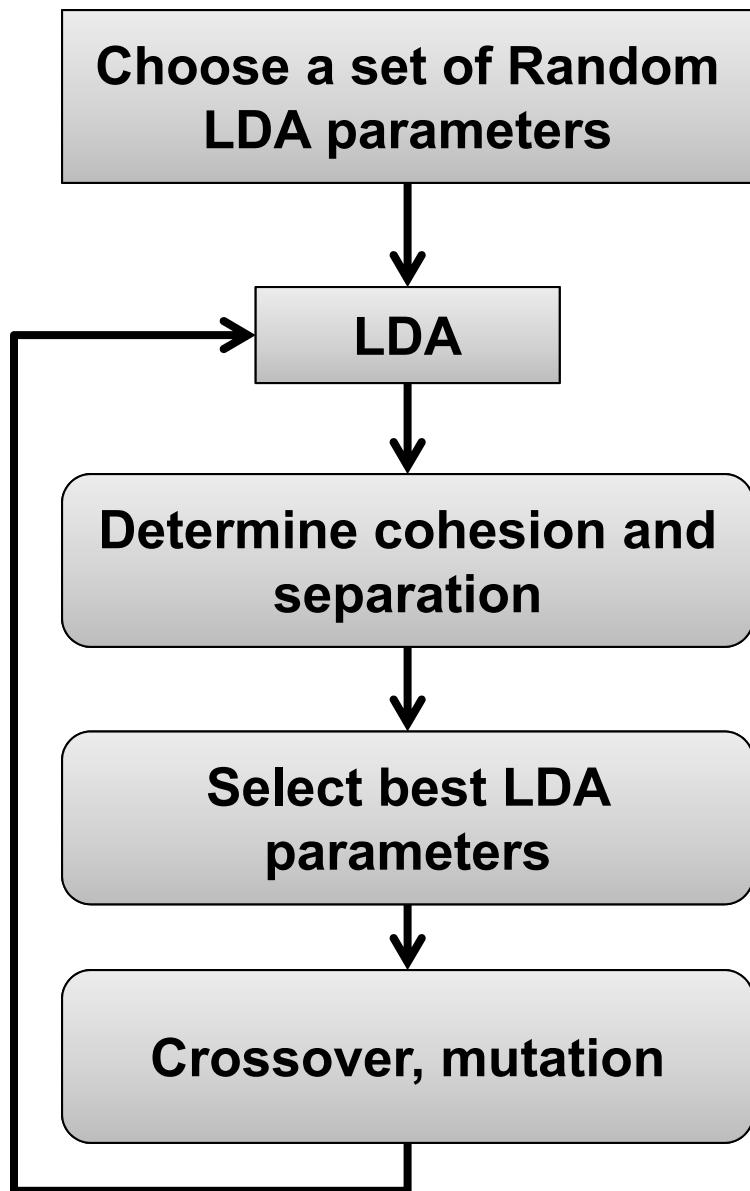


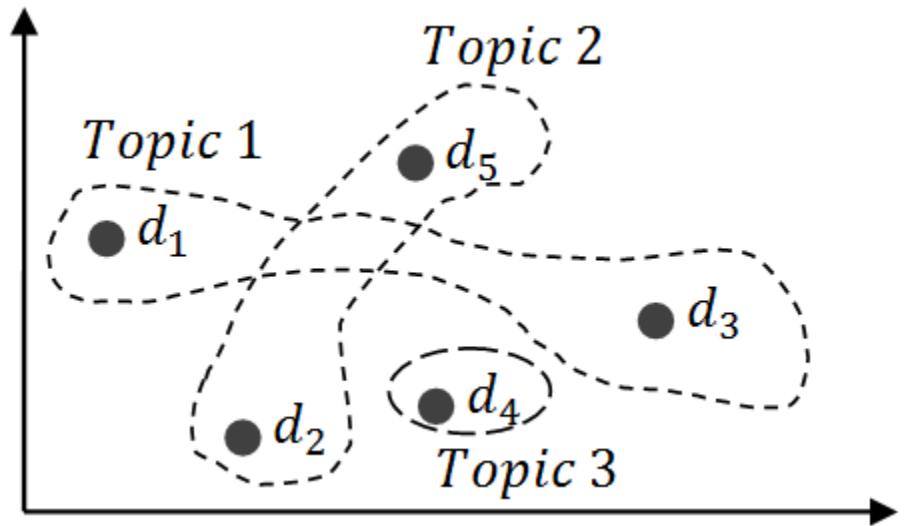
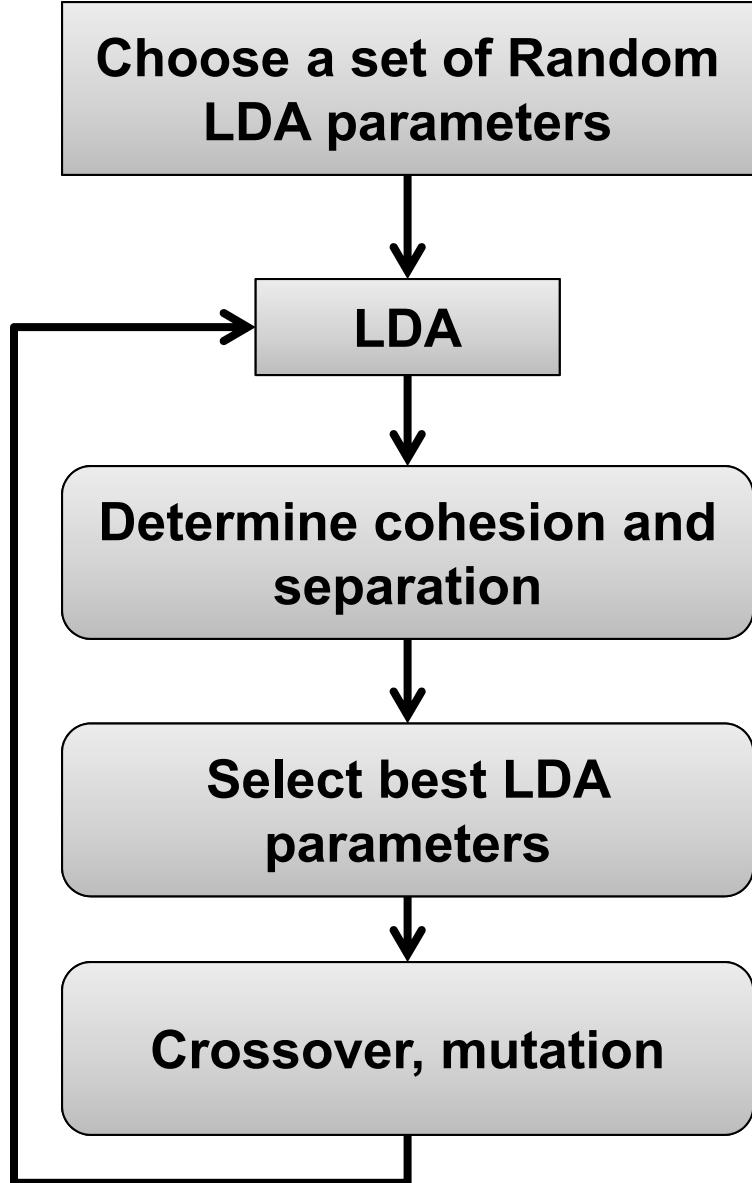
Crossover				
	iteration	topics	α	β
LDA Cfg. i	510	74	0.34	2.5
LDA Cfg. j	725	128	1.28	0.4
LDA Cfg. ij				

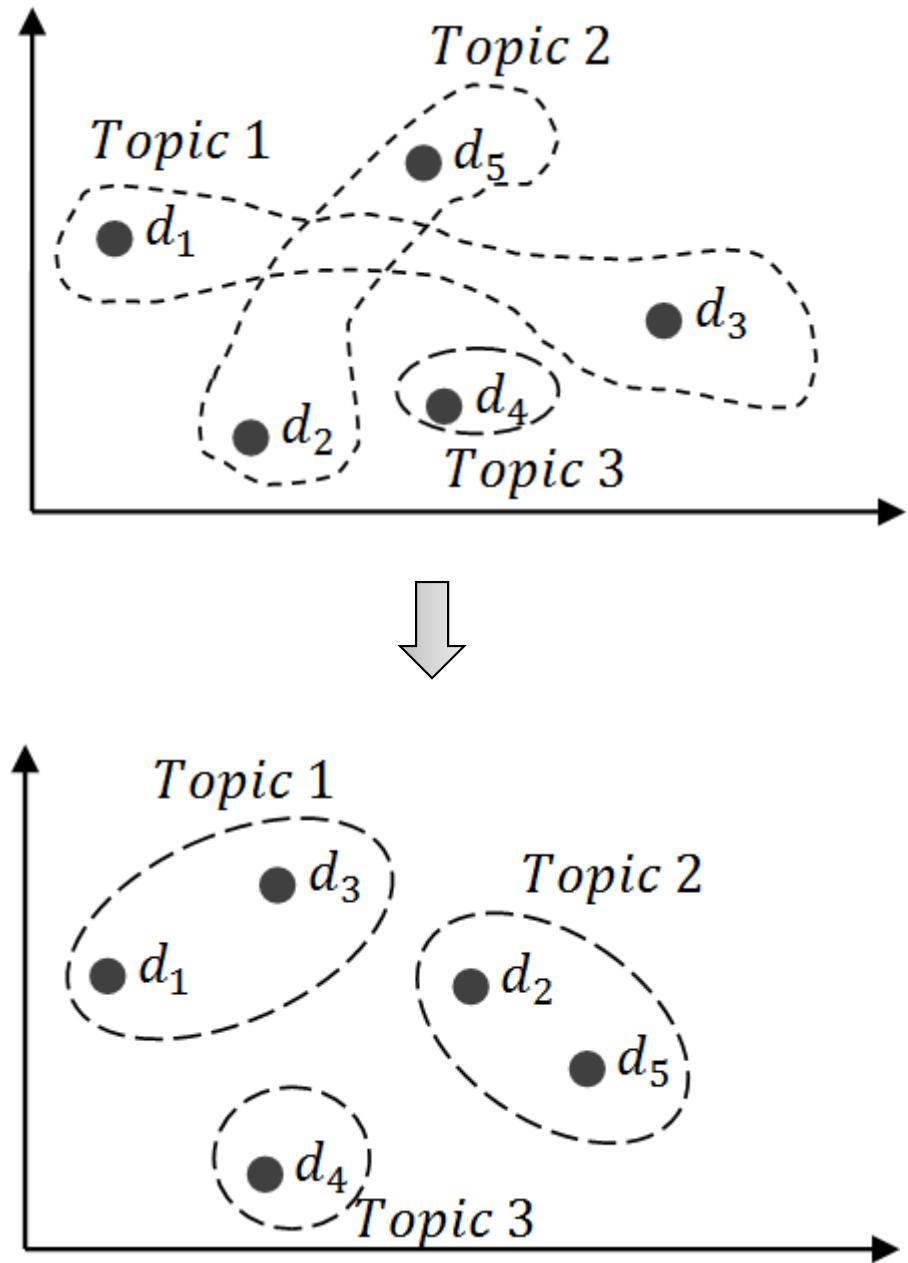
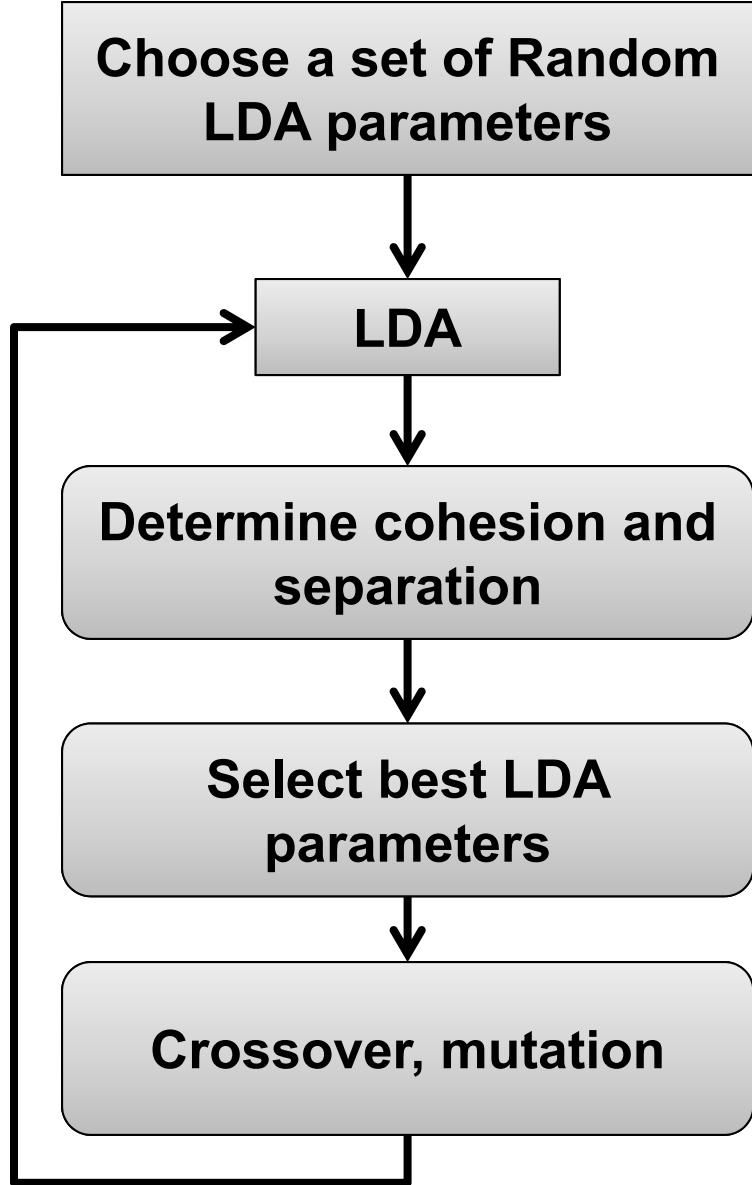


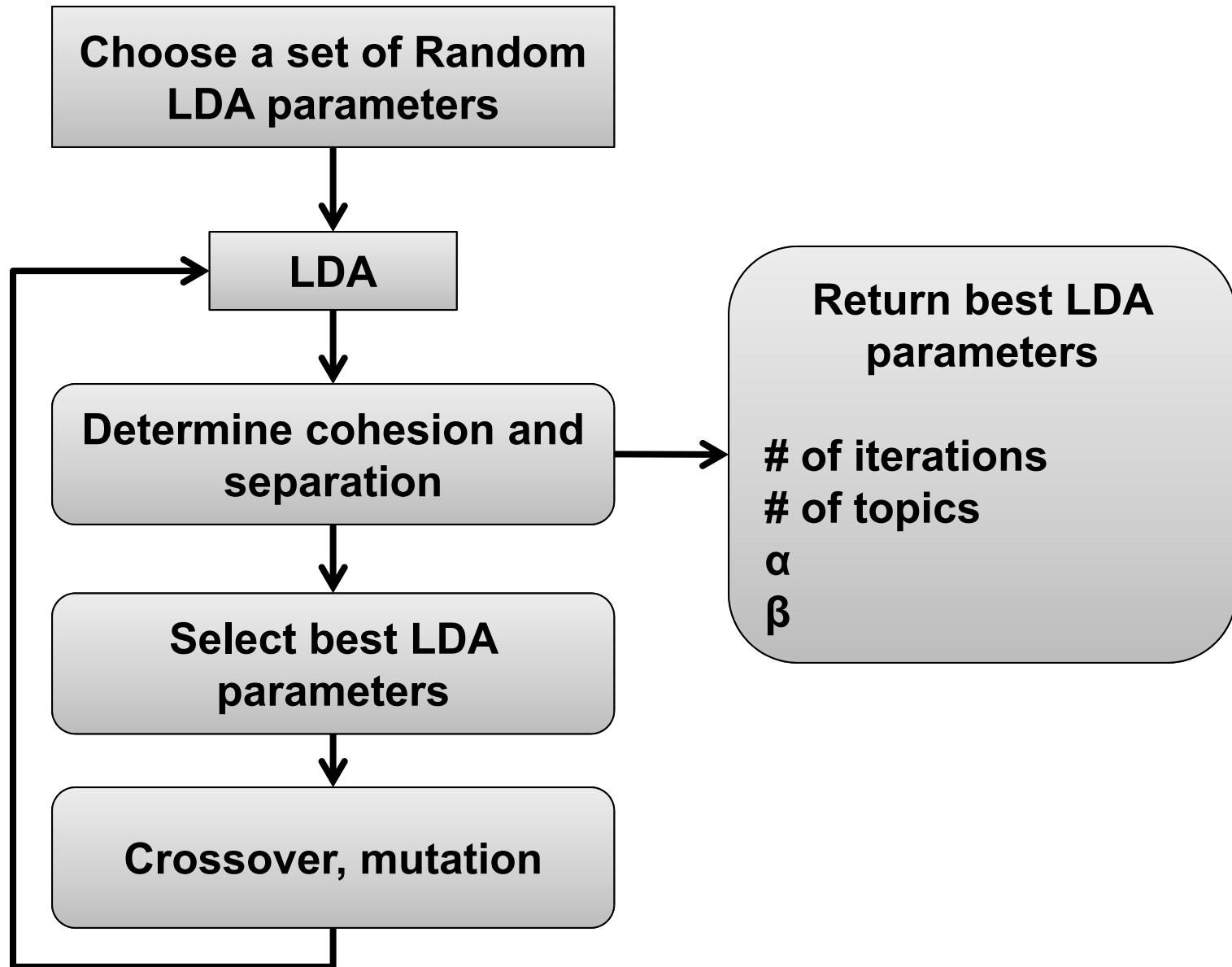




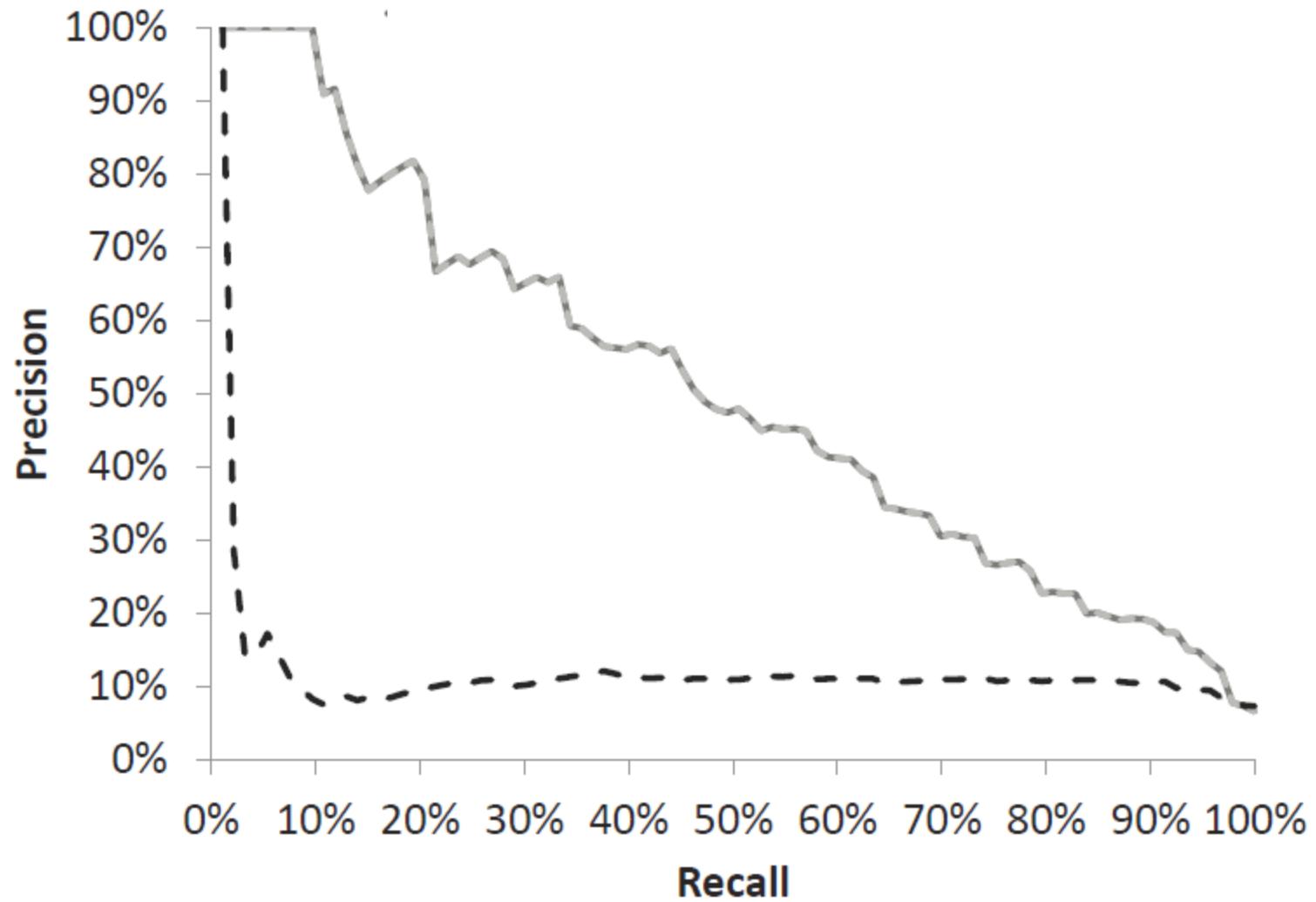




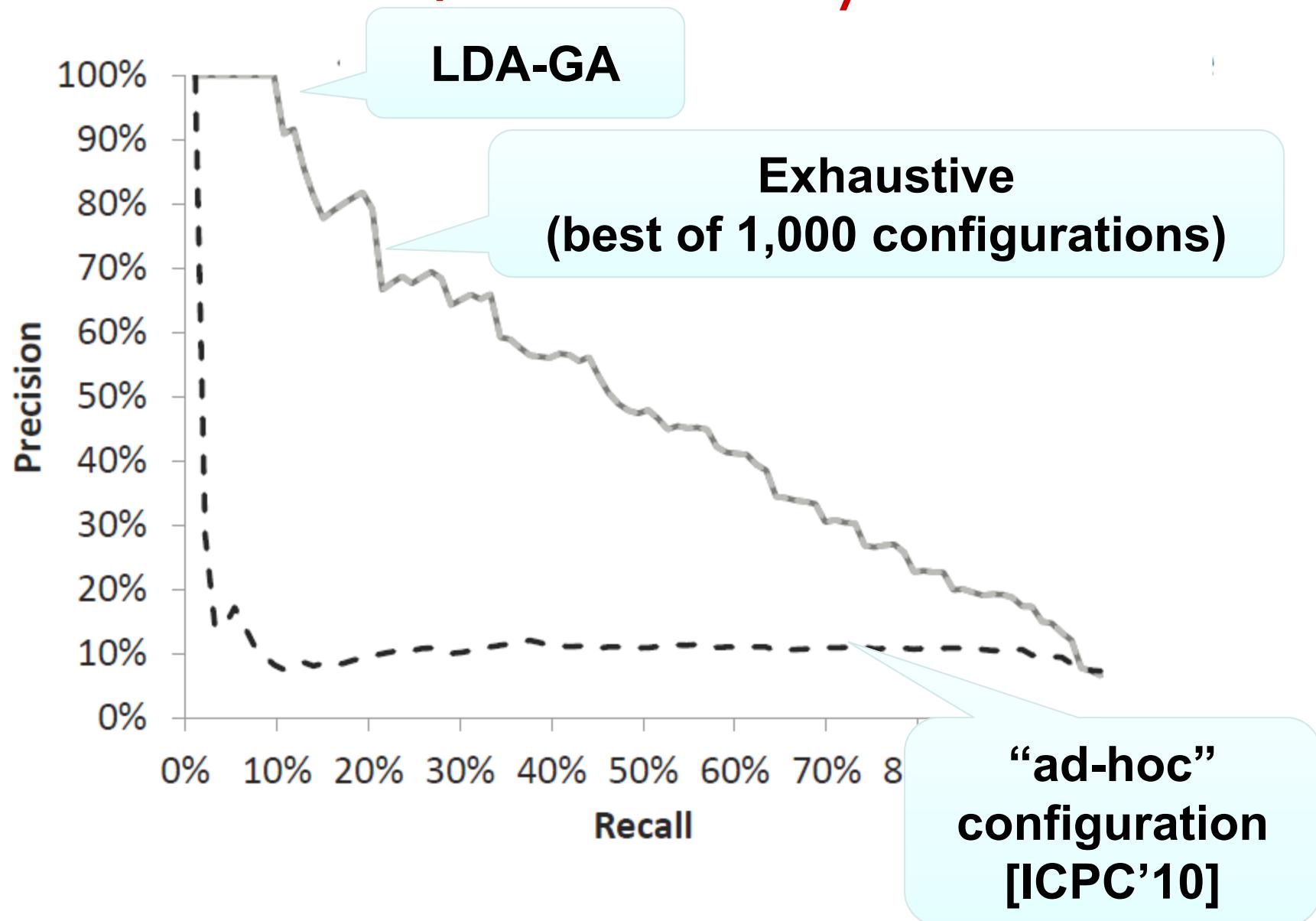




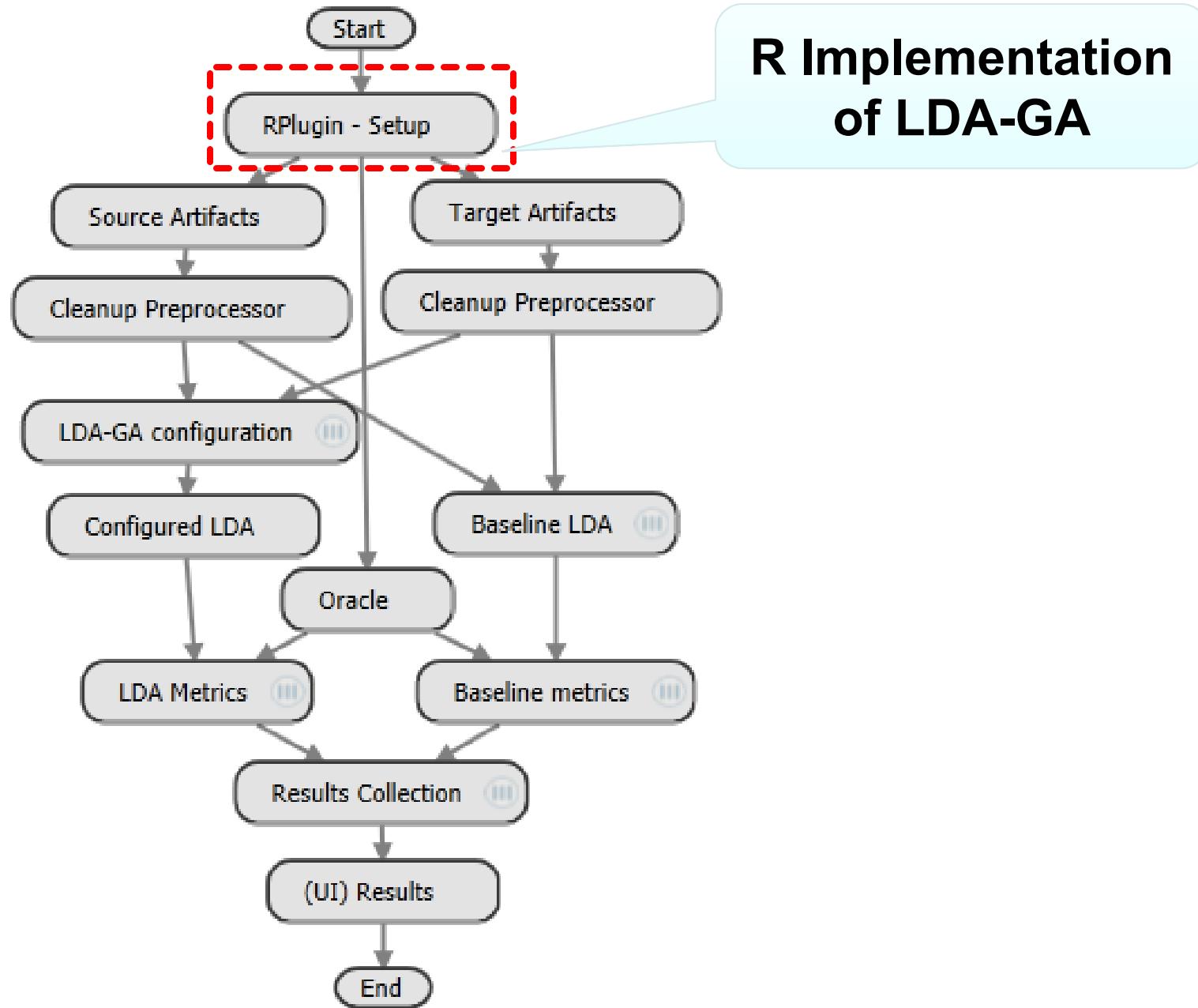
Precision/Recall EasyClinic



Precision/Recall EasyClinic

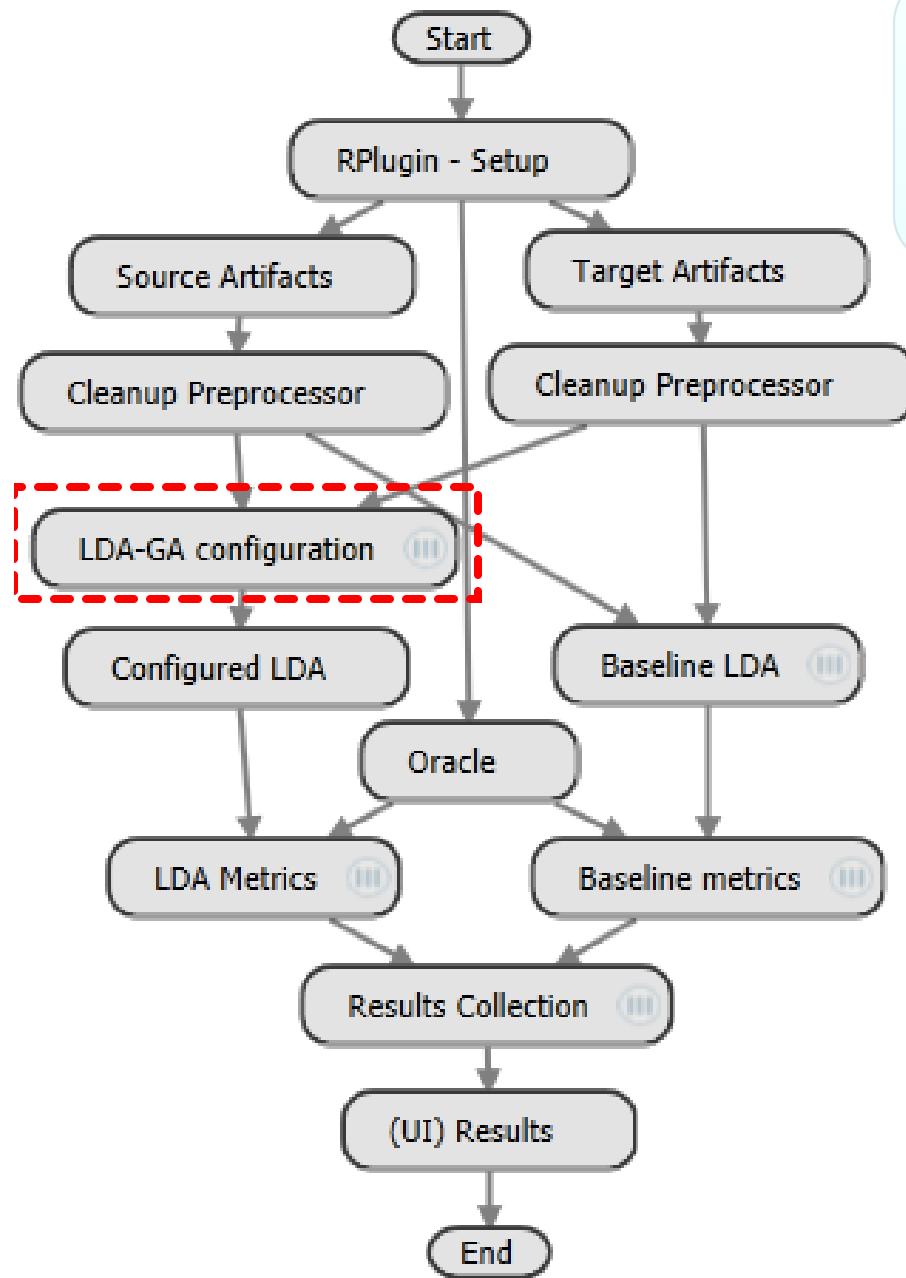


LDA-GA in TraceLab



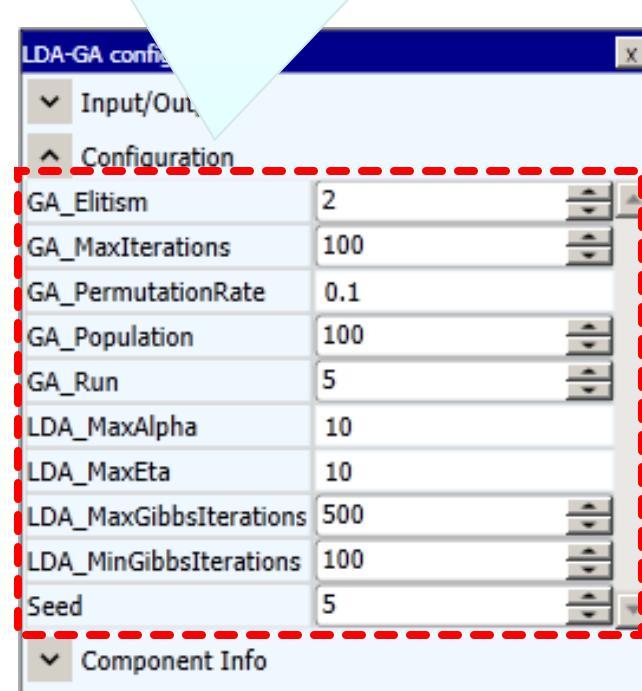
R Implementation
of LDA-GA

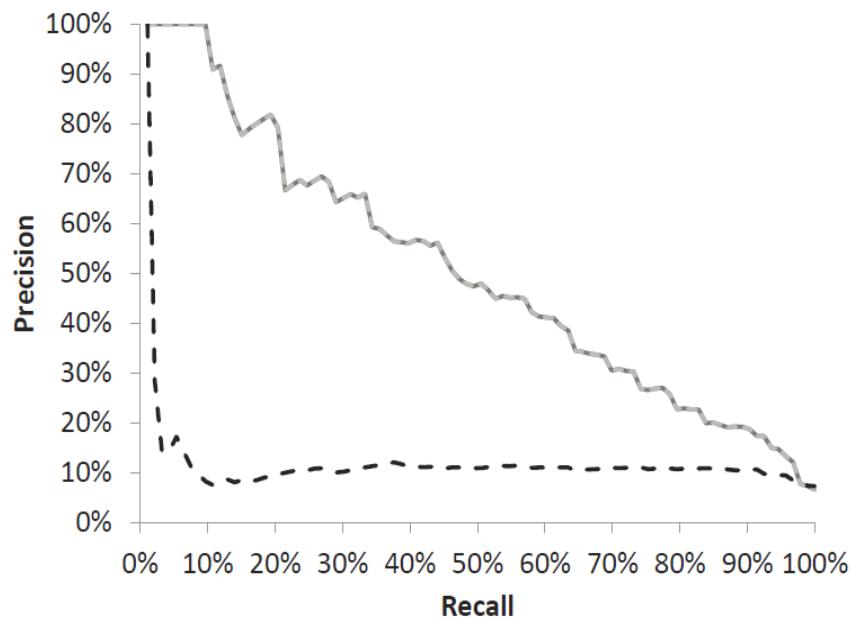
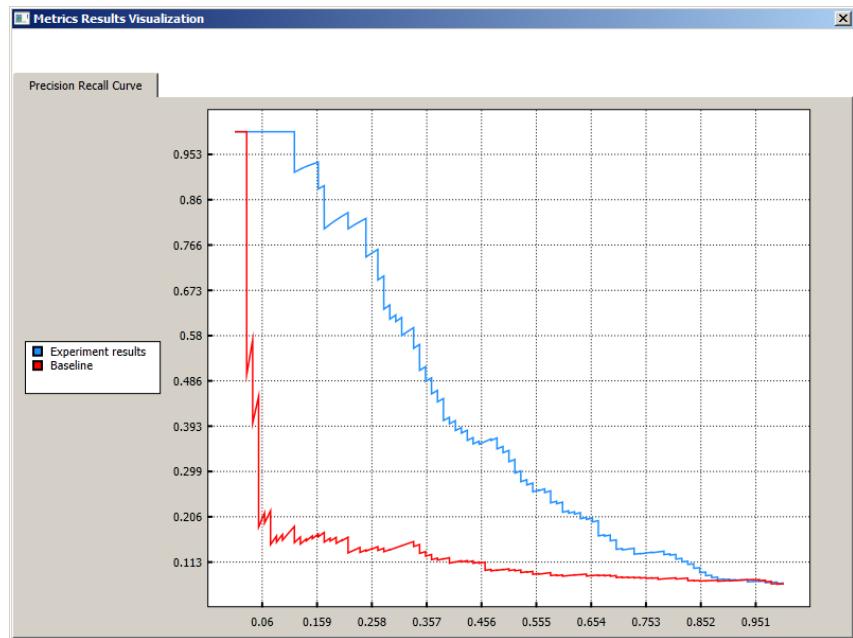
LDA-GA in TraceLab

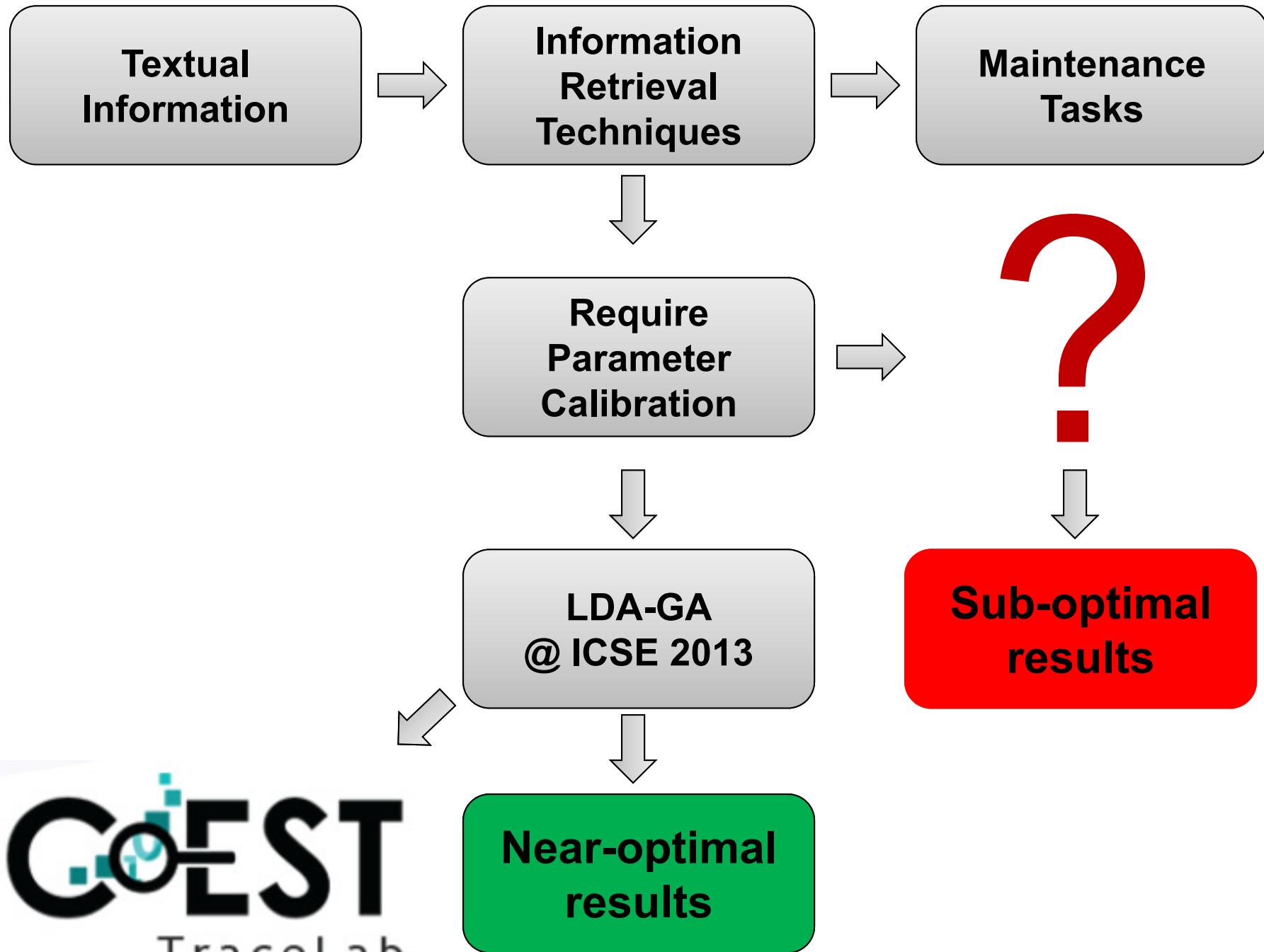


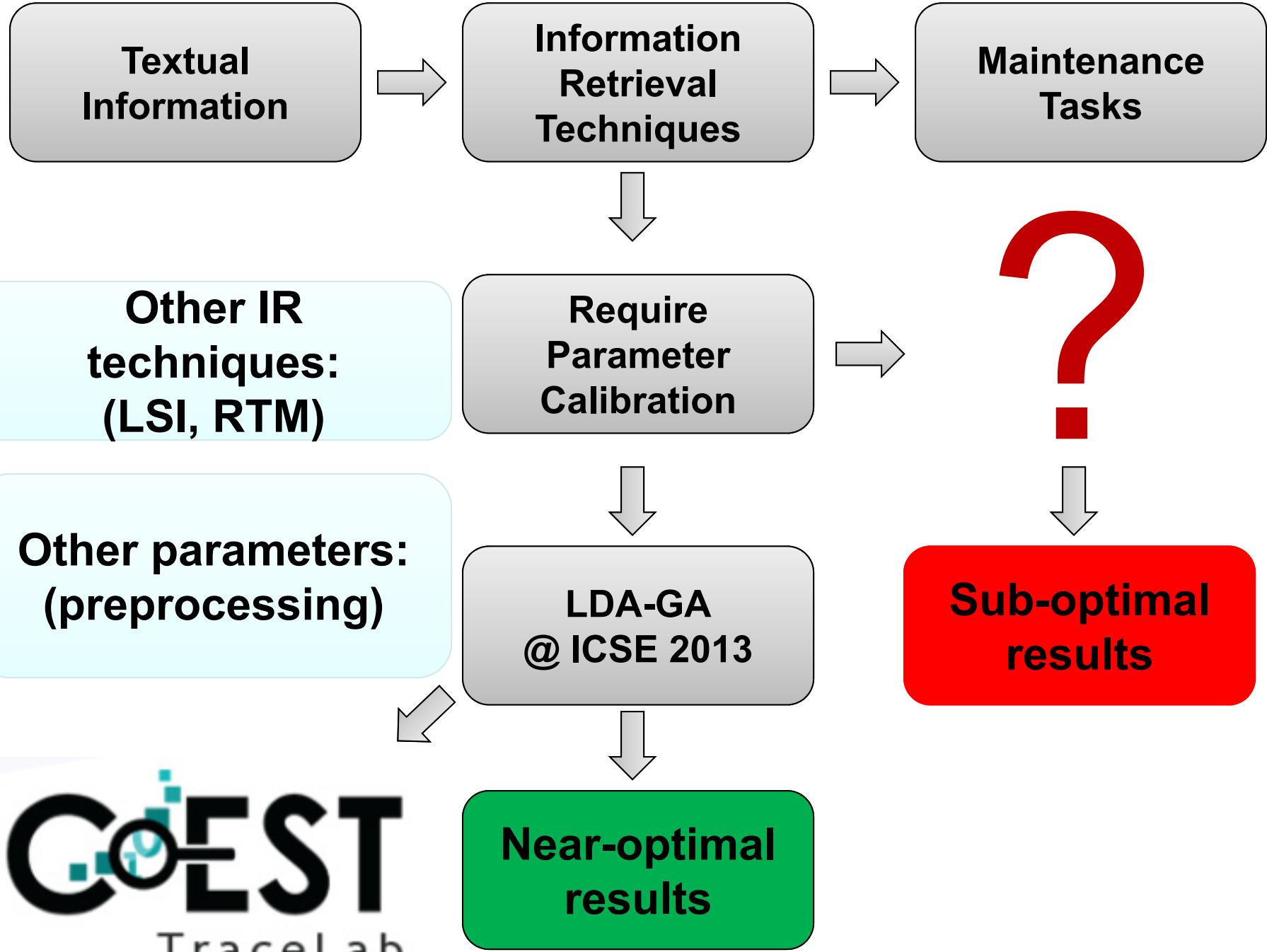
Genetic algorithm settings

LDA-GA settings





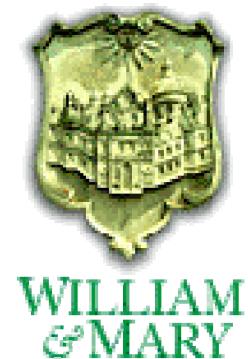




Thank you! Questions?

<http://www.cs.wm.edu/semeru/data/tefse13/>

<http://www.distat.unimol.it/reports/LDA-GA/>



UNIVERSITÀ
DEGLI STUDI
DEL MOLISE

