

Matrix and Tensor Factorization with Scientific Constraints

Ramakrishnan (Ramki) Kannan

Acknowledgements

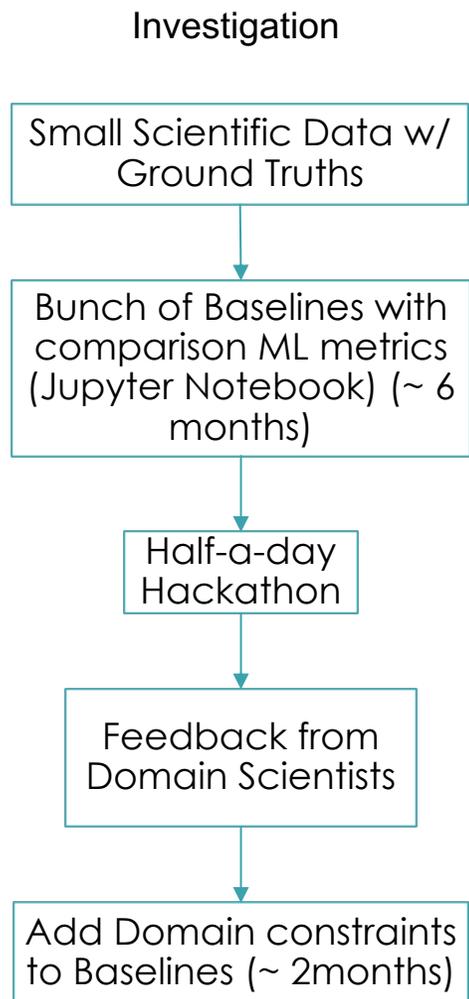
This manuscript has been co-authored by UT-Battelle, LLC under Contract No. DE-AC05-00OR22725 with the U.S. Department of Energy. This project was partially funded by the Laboratory Director's Research and Development fund. This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy.

This research used resources of the National Energy Research Scientific Computing Center, a DOE Office of Science User Facility supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231

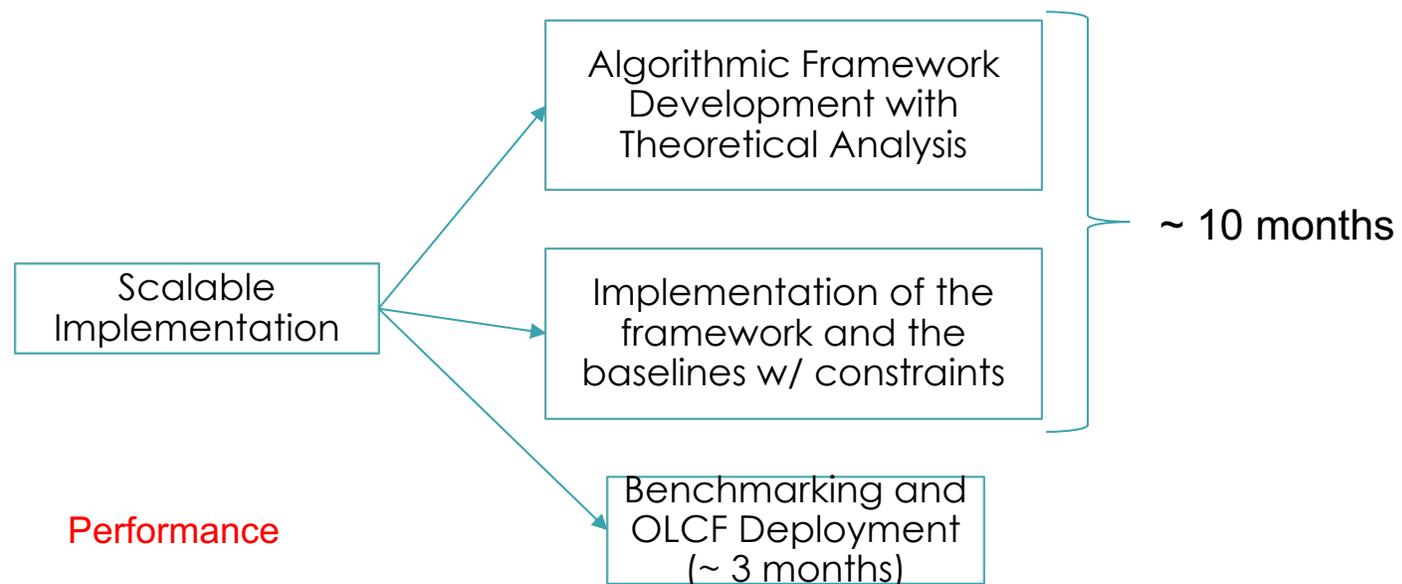
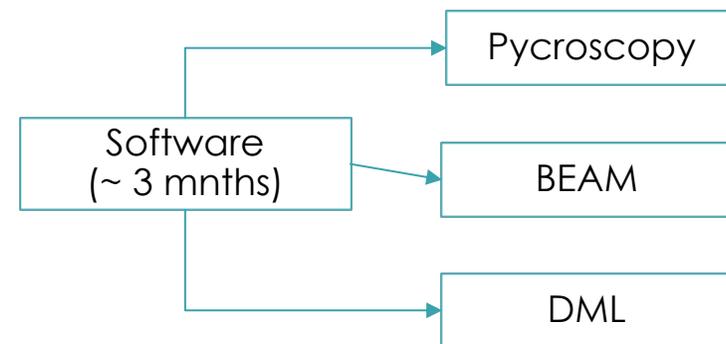
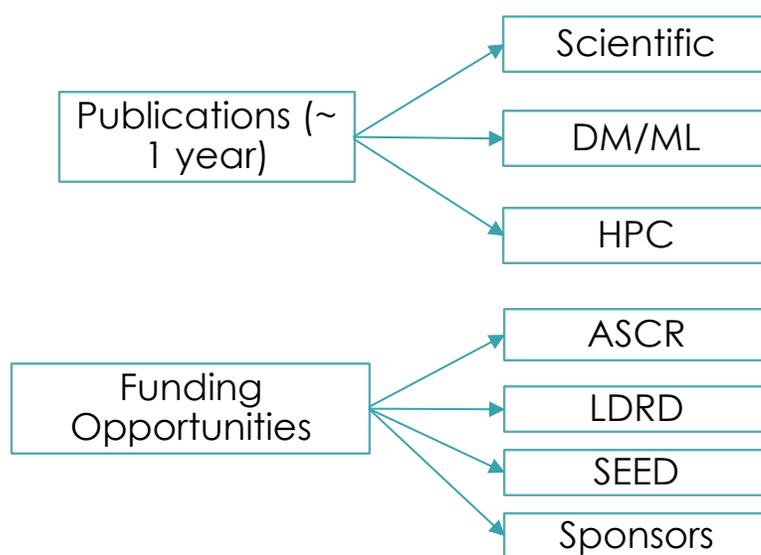
The United States Government retains and the publisher, by accepting the article for publication, acknowledges that the United States Government retains a non-exclusive, paid-up, irrevocable, world-wide license to publish or reproduce the published form of this manuscript, or allow others to do so, for United States Government purposes. The Department of Energy will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan

<http://energy.gov/downloads/doepublic-access-plan>. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the USDOE.

Scientific Engagement Model

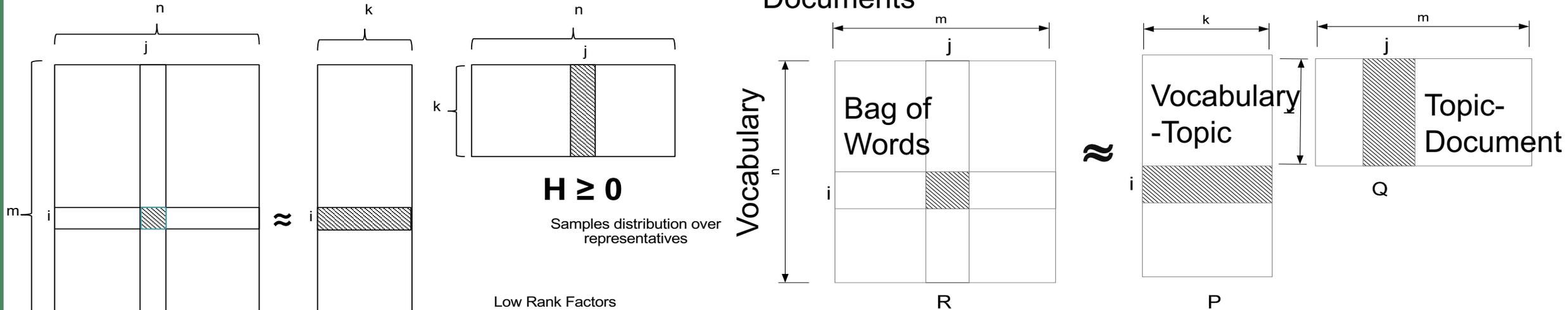


Productivity



Performance

NMF and Applications



Low Rank Factors

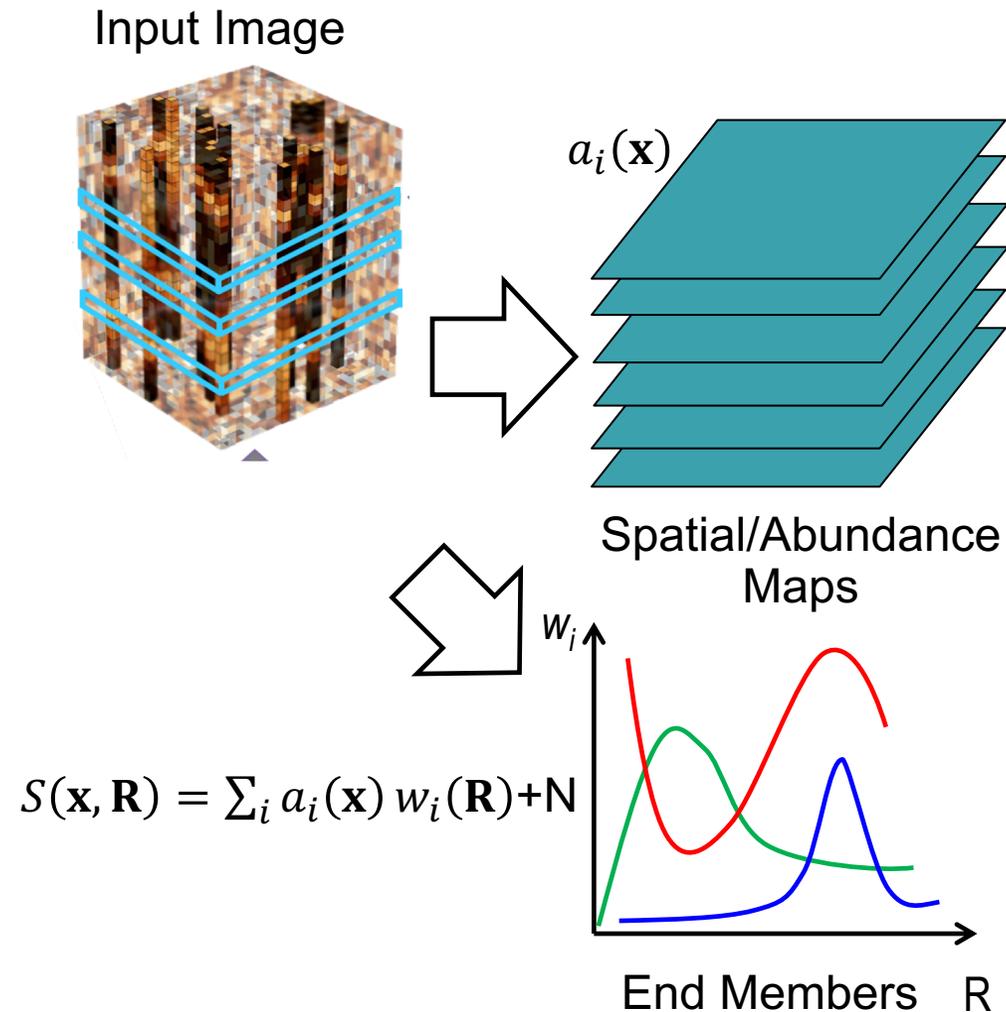
Top Keywords from Topics 1-25					Top Keywords from Topics 26-50				
word1	word2	word3	word4	word5	word1	word2	word3	word4	word5
refer	undefin	const	key	compil	echo	type=text	php	form	result
text	field	box	word	static	test	perform	fail	unit	result
imag	src	descript	alt=ent	size	tabl	key	queri	databas	insert
button	click	event	form	add	user	email	usernam	login	log
creat	bean	add	databas	except	data	json	store	read	databas
string	static	final	catch	url	page	load	content	url	link
width	height	color	left	display	privat	static	final	import	float
app	applic	servic	thread	work	row	column	date	cell	valu
ipsum	lorem	dolor	sit	amet	line	import	command	print	recent
node	list	root	err	element	var	map	marker	match	url
0x00	0xff	byte	0x01	0xc0	server	connect	client	messag	request
file	directori	read	open	upload	number	byte	size	print	input

Ramakrishnan Kannan, [Grey Ballard](#), [Haesun Park](#): MPI-FAUN: An MPI-Based Framework for Alternating-Updating Nonnegative Matrix Factorization. [IEEE Trans. Knowl. Data](#)

Enq.30(3): 544-558 (2018)

Motivation

- Understanding terrestrial information in an unknown place from satellite images
- Identifying presence of hidden unknown/foreign bodies in a scanned image - Eg., contamination in food articles, camouflaged explosives etc.
- Biological application - spectral karyotyping, immunofluorescence, live-cell imaging, drug discovery, and tissue pathology – Eg., Unmixing on Spectral imaging of the stained tissues using multiple dyes.
- Physics and Material Sciences – Mapping properties to end-members. Comparing different materials



MPI-FAUN

- Distributed Communication avoiding NMF Algorithms
- <https://github.com/ramkikannan/nmflibrary>
- <http://dx.doi.org/10.1109/TKDE.2017.2767592>

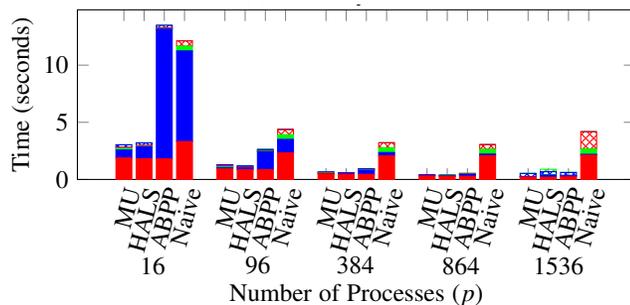
Titan – Dense Matrix, Low Rank 50, 100 Iterations, 12650 Nodes, 202500 Cores,

Matrix Size	Algos	NMF Time (in Secs)
3.03 million x 3.03 million	MU	554
	HALS	197.75
	ANLS/BPP	219.8

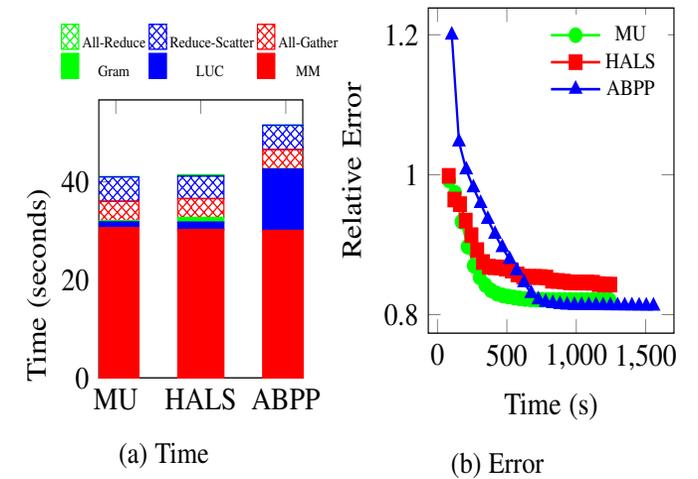
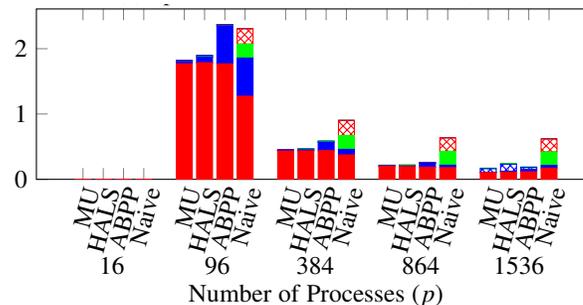
Rhea, 100 nodes, 1600 cores, Low Rank 50,

Dataset	Type	Matrix size	NMF Time
Video	Dense	1 Million x 13,824	5.73 seconds
Stack Exchange	Sparse	627,047 x 12 Million	67 seconds
Webbase-2001	Sparse	118 Million x 118 Million	25 minutes

Sparse Webbase – 1 Million Vertex Graph



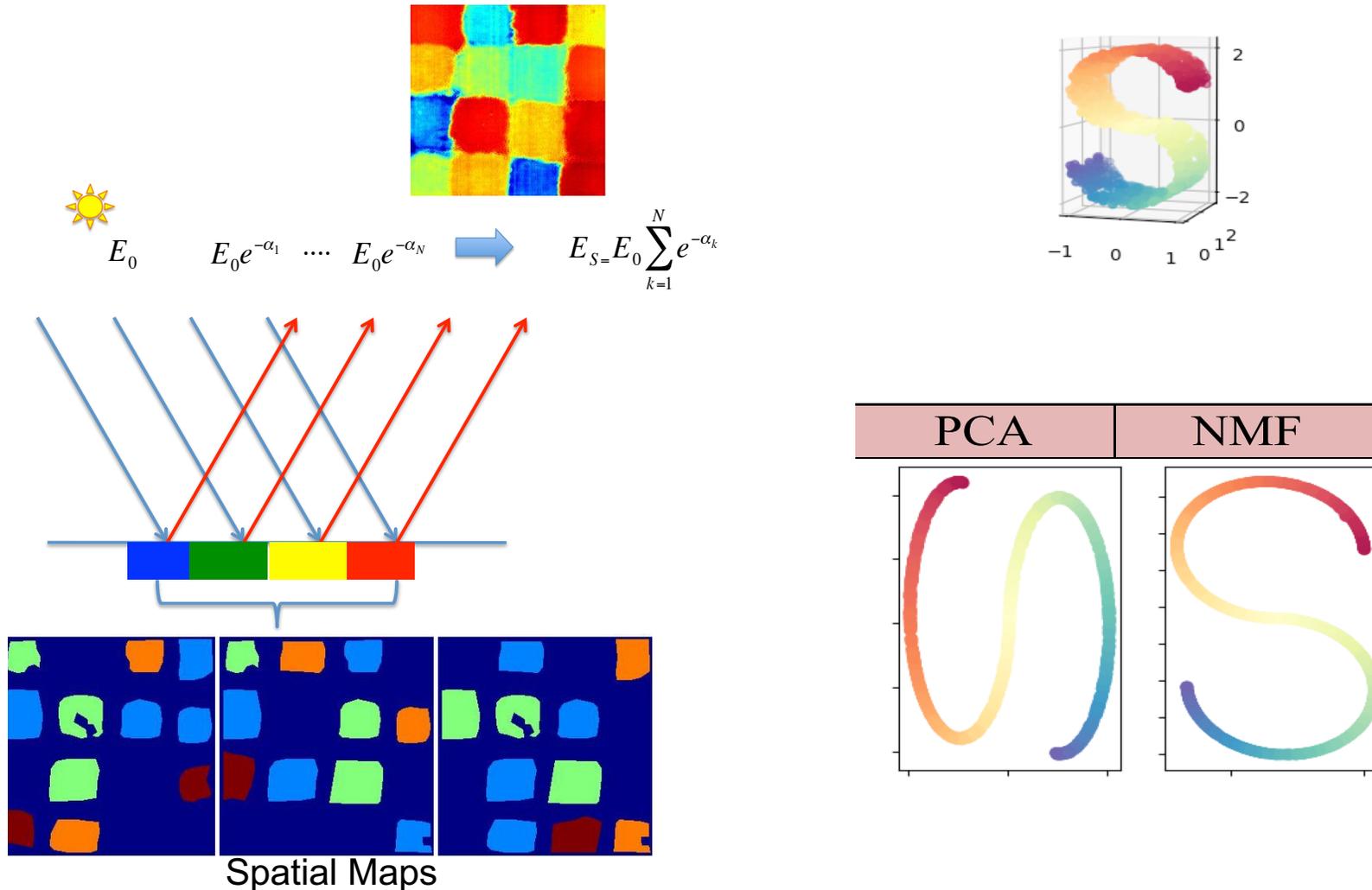
Dense Real world – Video



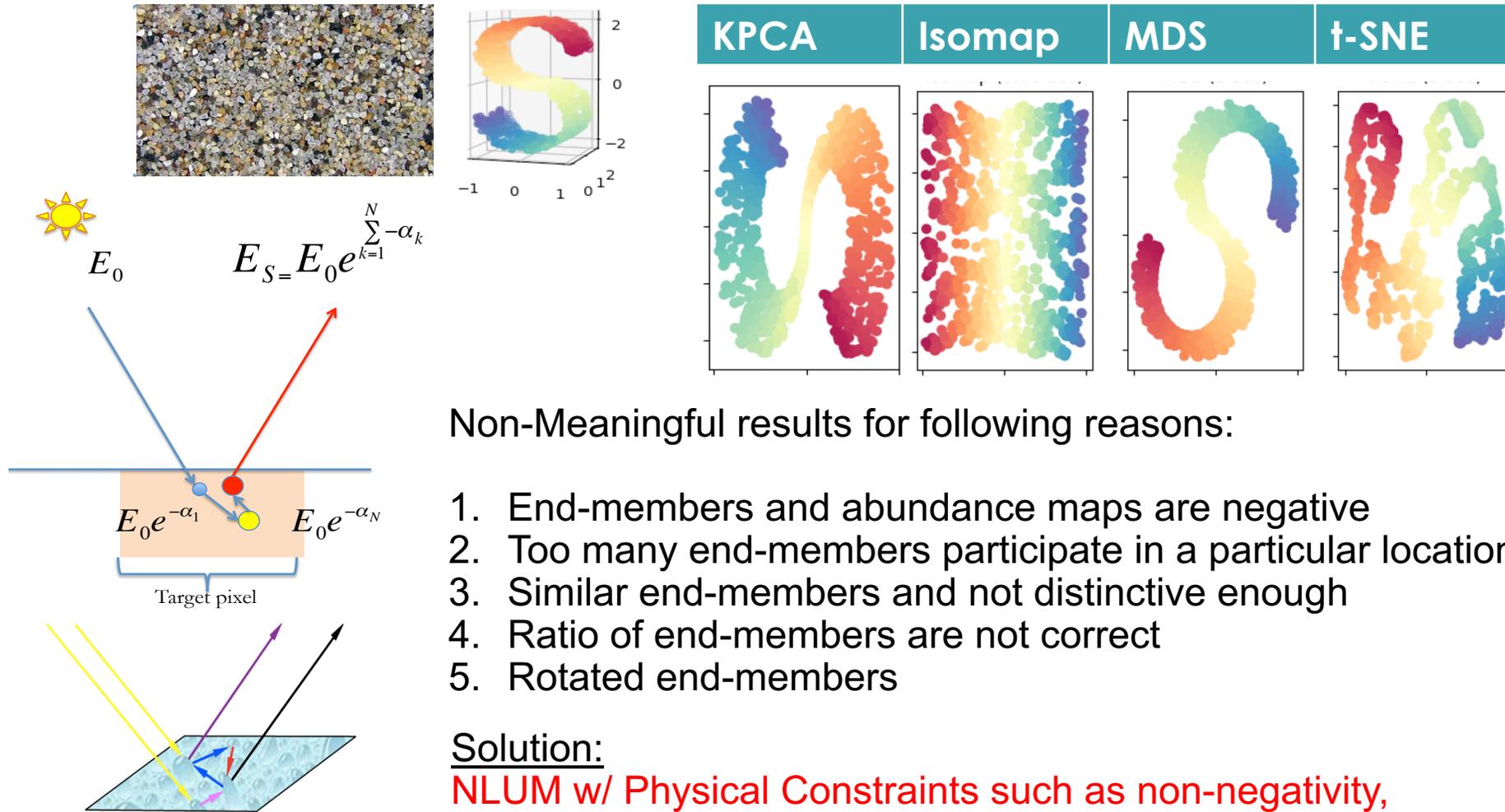
NMF on 118 million Web-graph

Existing Approach : Linear Unmixing

1. Good at Capturing Macroscopic Information
2. Spatially segregated patterns



Existing Non-linear Unmixing (NLUM)



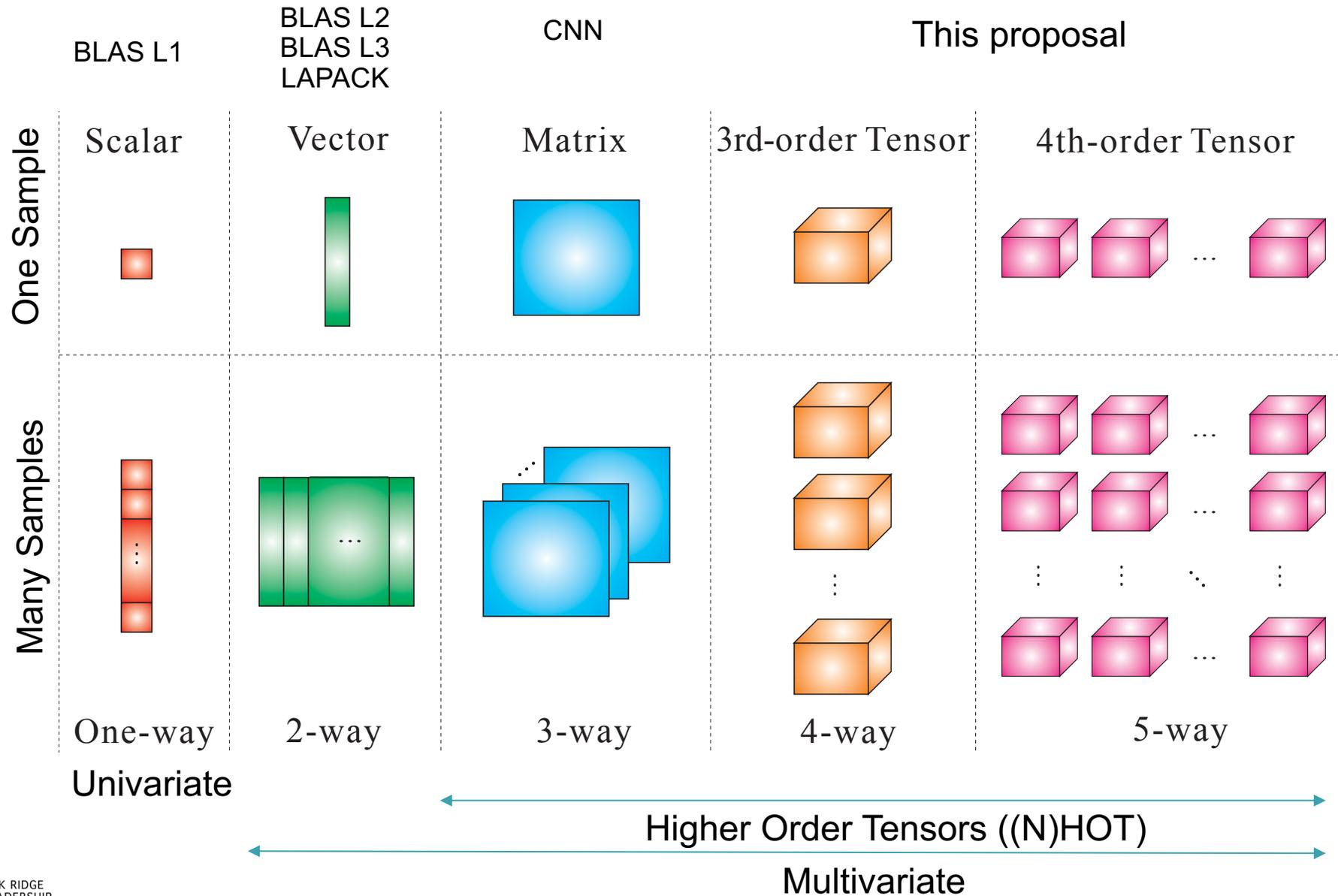
Non-Meaningful results for following reasons:

1. End-members and abundance maps are negative
2. Too many end-members participate in a particular location
3. Similar end-members and not distinctive enough
4. Ratio of end-members are not correct
5. Rotated end-members

Solution:

NLUM w/ Physical Constraints such as non-negativity, sparsity, spatial smoothness, sum to 1, orthogonal etc.

Higher Order Tensors



Dimensionality Reduction in Scientific Data

- Multimodal characterization of materials – *comprehensive characterization from chemical composition to functional properties on the nanoscale*

