

Identify Dwarfs Workloads in Big Data Analytics

Xinhui Tian

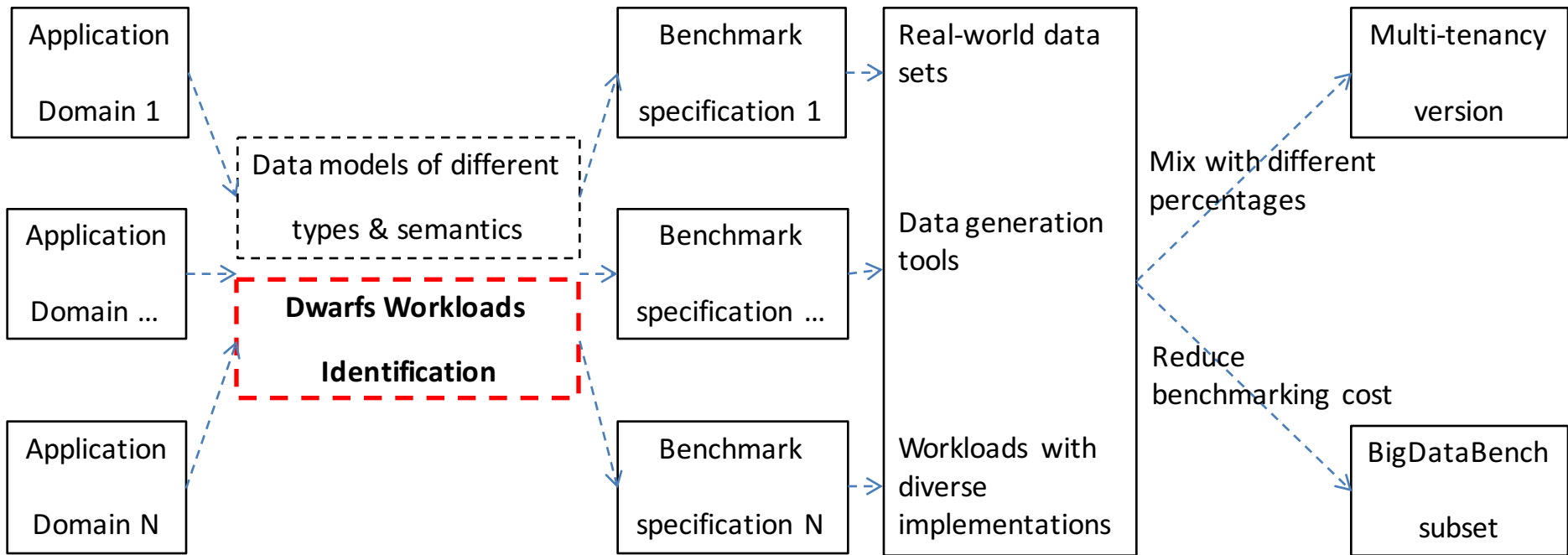
***Institute of Computing Technology,
Chinese Academy of Sciences***

**BigDataBench Tutorial
ASPLOS 2016 Atlanta, GA**



中国科学院
INSTITUTE OF COMPUTING TECHNOLOGY

BigDataBench Methodology



How to define a representative big data benchmark ?

- One attempt
 - Using
 - Subject
- Another attempt
 - Gains of

How can we construct a benchmark suite using **a minimum workload set** to represent **maximum patterns** of big data analytics?

Dwarfs workloads !



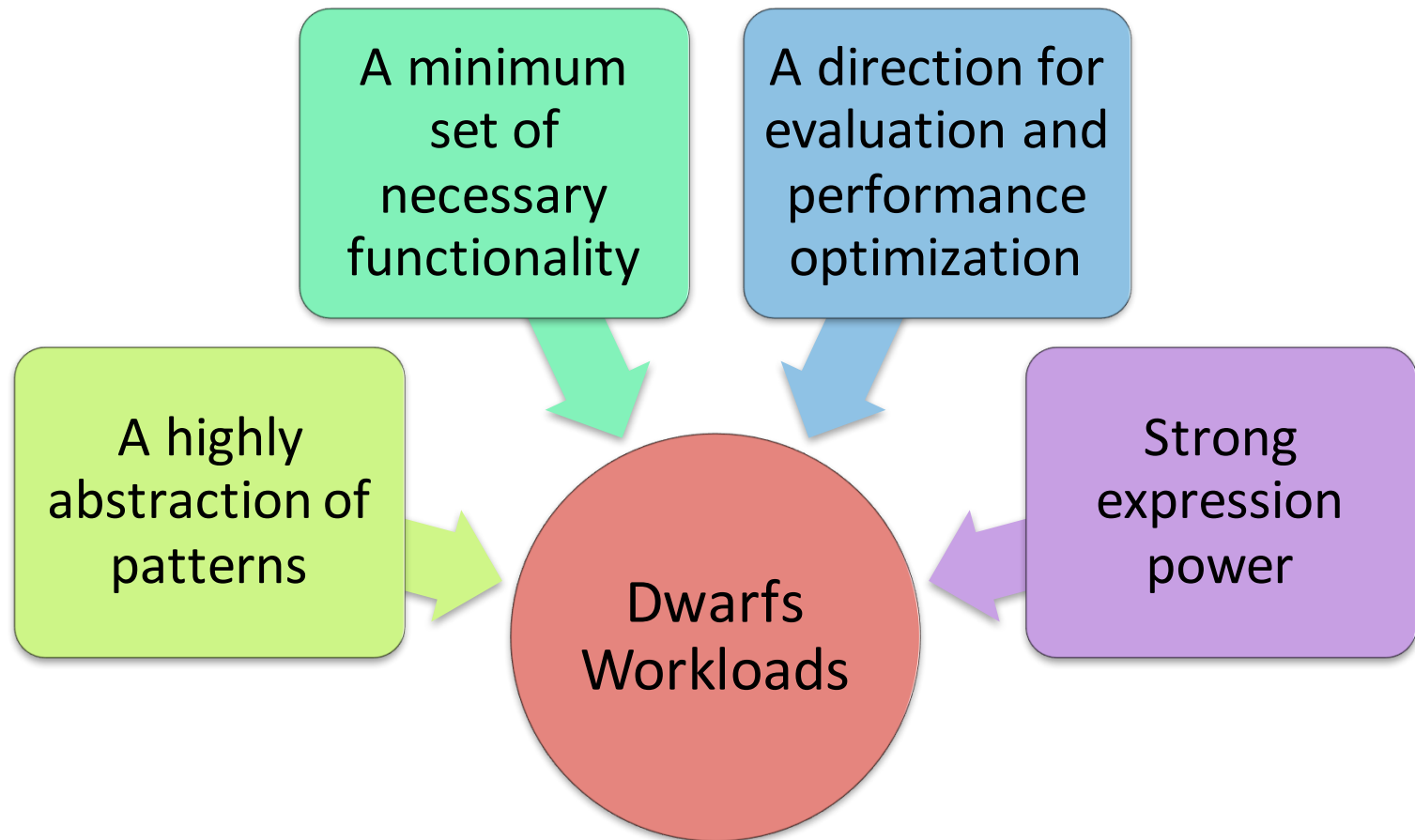
Inspiration

Successful Compute Abstractions *Successful Benchmarks*

- Relational algebra
 - 5 primitive operations
 - Select, Project, Product, Union, Difference
- Parallel computing
 - Computational & communication patterns
 - 13 dwarfs
- TPC-C
 - OLTP domain
 - Functions of abstraction
- HPCC
 - High performance computing
 - Seven basically tests

K. Asanovic,, etc. The landscape of parallel computing research: A view from berkeley.

Why Dwarfs are Important



<http://www.krellinst.org/doecsgf/conf/2014/pres/jhill.pdf>

<http://cacs.usc.edu/education/cs596/DavidPatterson.pdf>

Fundamental Issues

What are the challenges for big data dwarfs?

How to find the dwarfs workloads in big data analytics?



Challenges

Structured Data



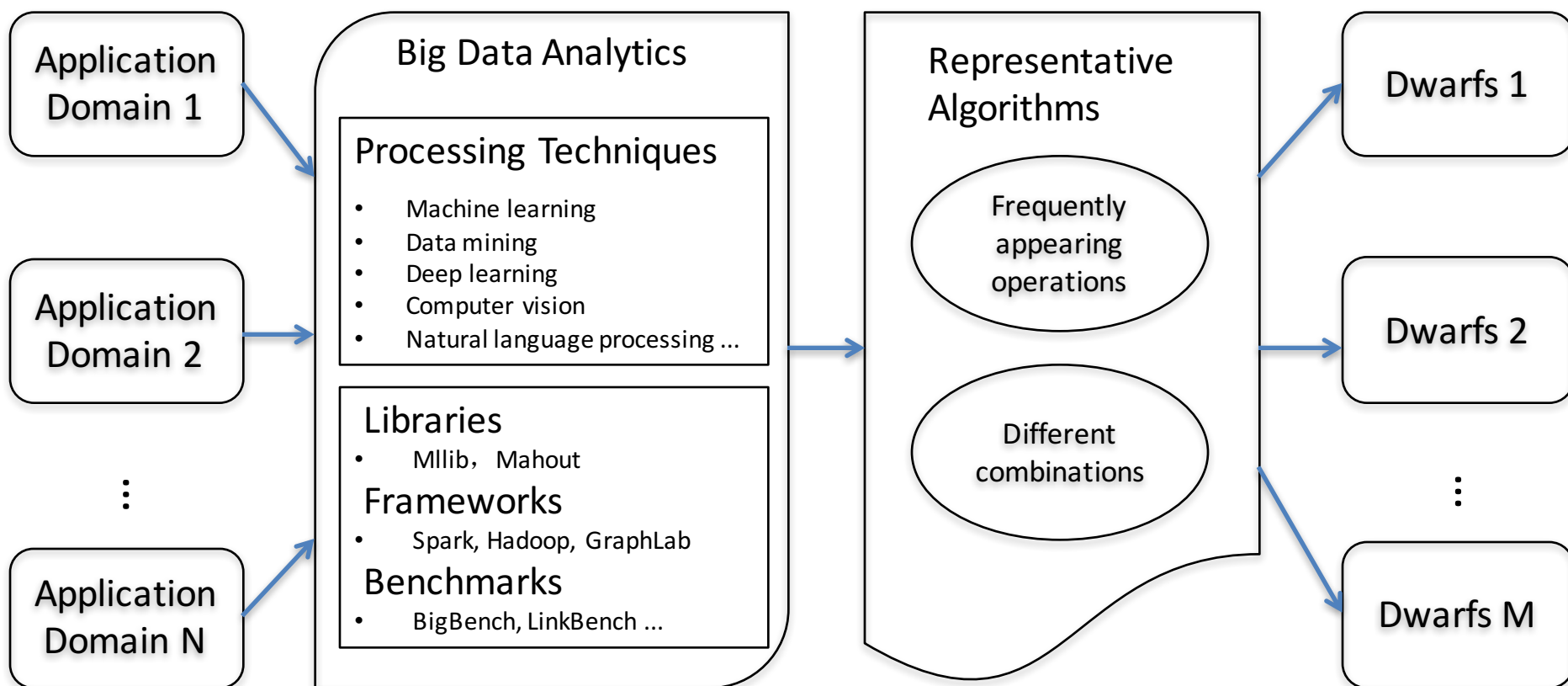
0.103	0.176	0.387	0.300	0.379
0.333	0.384	0.564	0.587	0.857
0.421	0.309	0.654	0.729	0.228
0.266	0.750	1.056	0.936	0.911
0.225	0.326	0.643	0.337	0.721
0.187	0.586	0.529	0.340	0.829
0.153	0.485	0.560	0.428	0.628

Unstructured Data



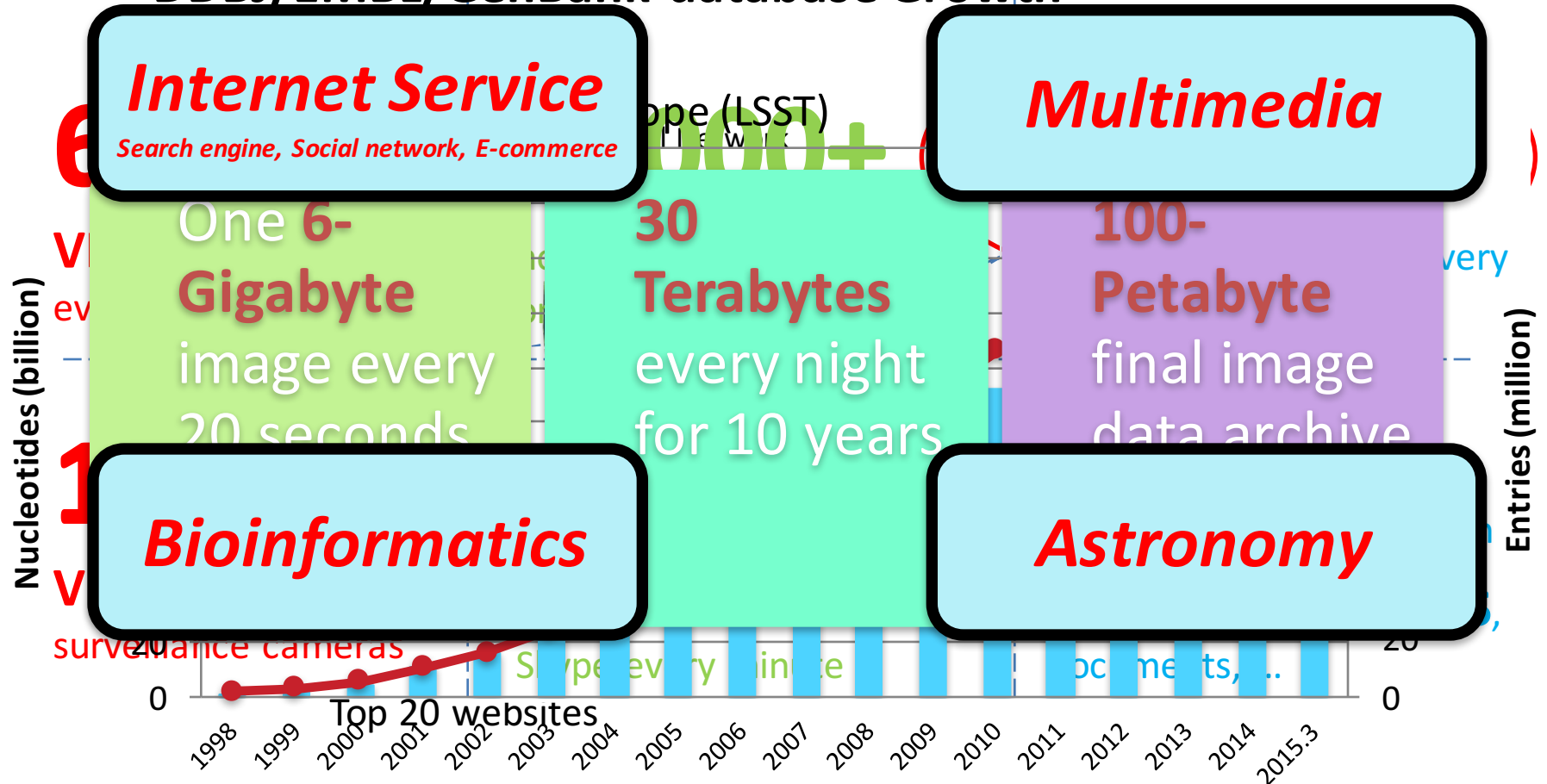
- *80%* data growth are unstructured data
- Operations on big data are *complicated*
 - Pipeline? Parallel?

Methodology of Dwarfs

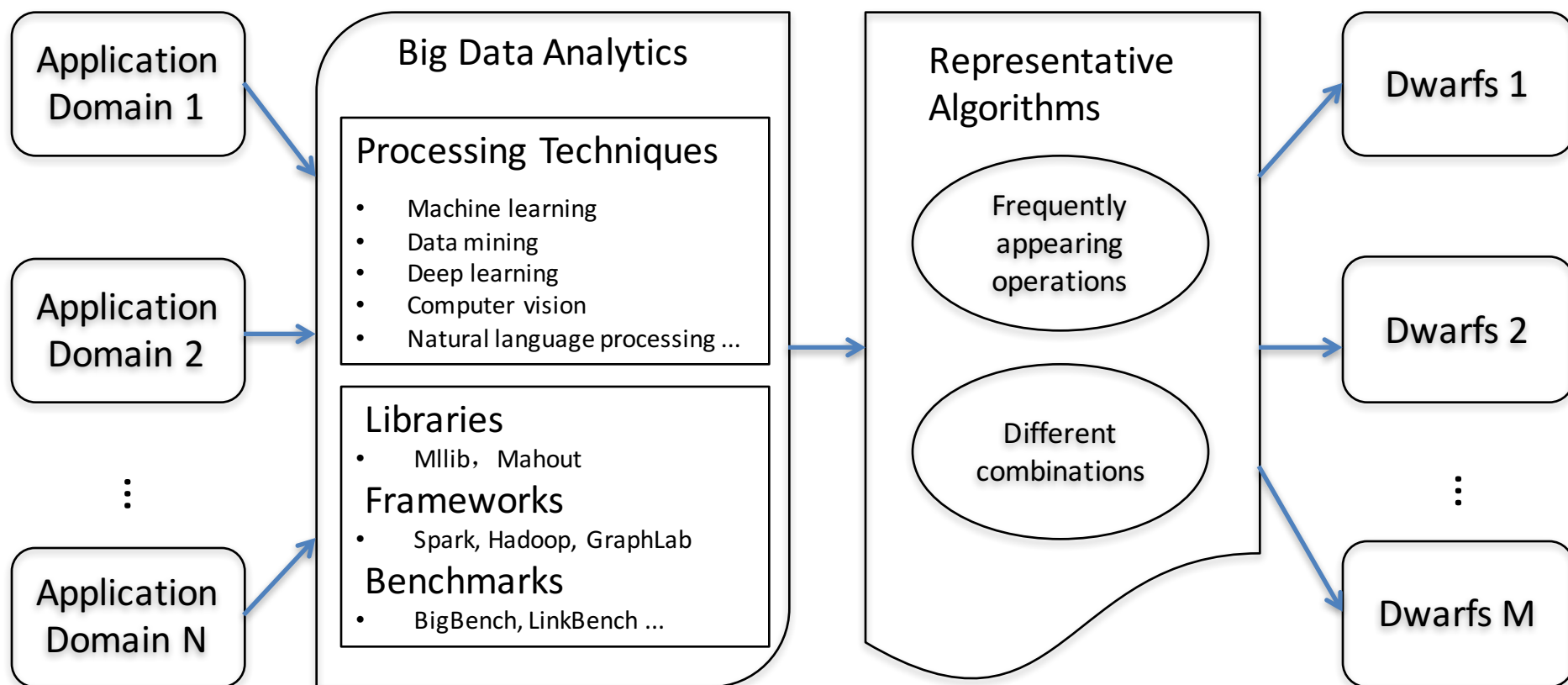


Choosing Application Domain

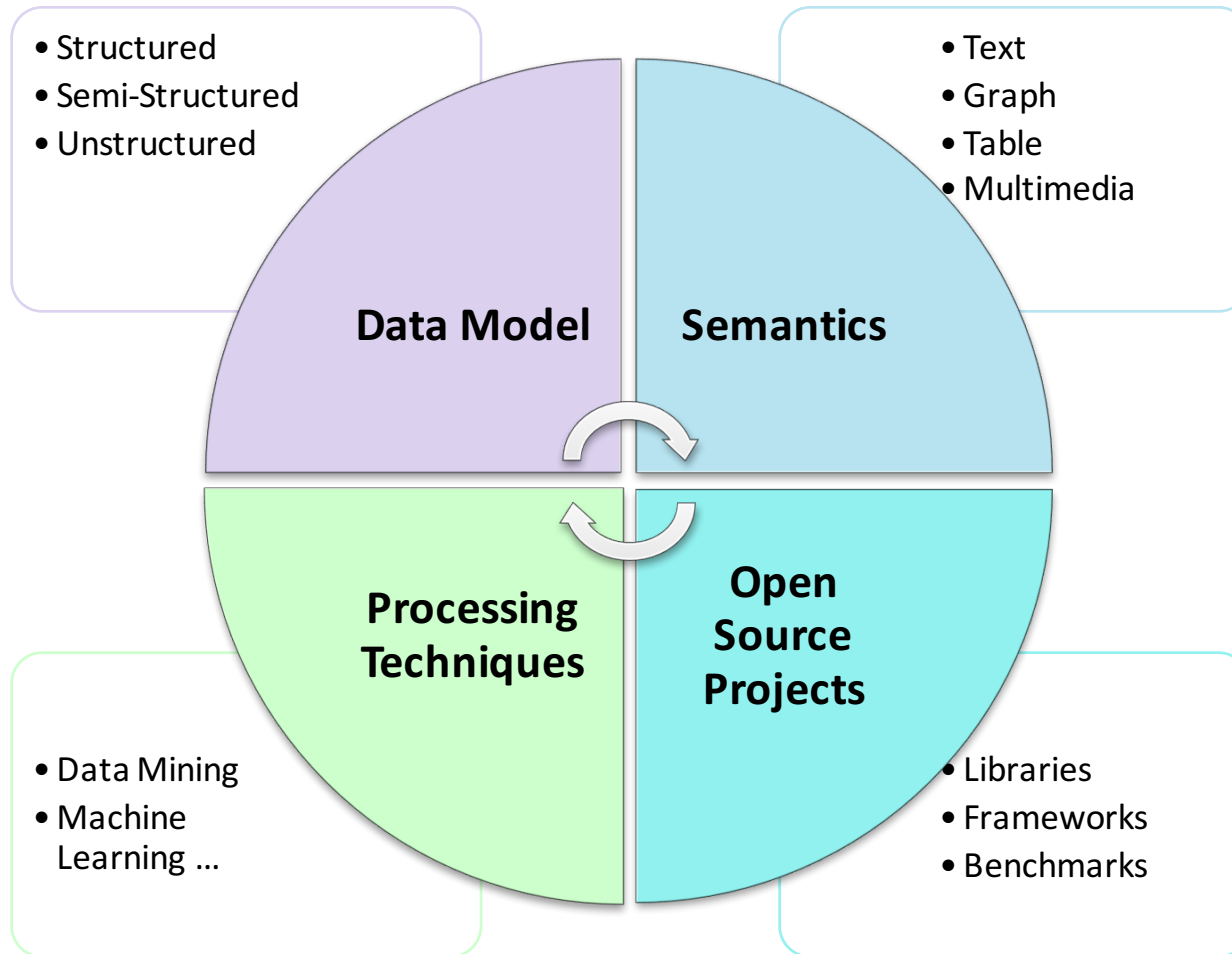
DDBJ/EMBL/GenBank database Growth



Methodology of Dwarfs



Big Data Analytics



Libraries & Frameworks & Benchmarks

Libraries

- Opencv, Mlib, Weka, AstroML ...

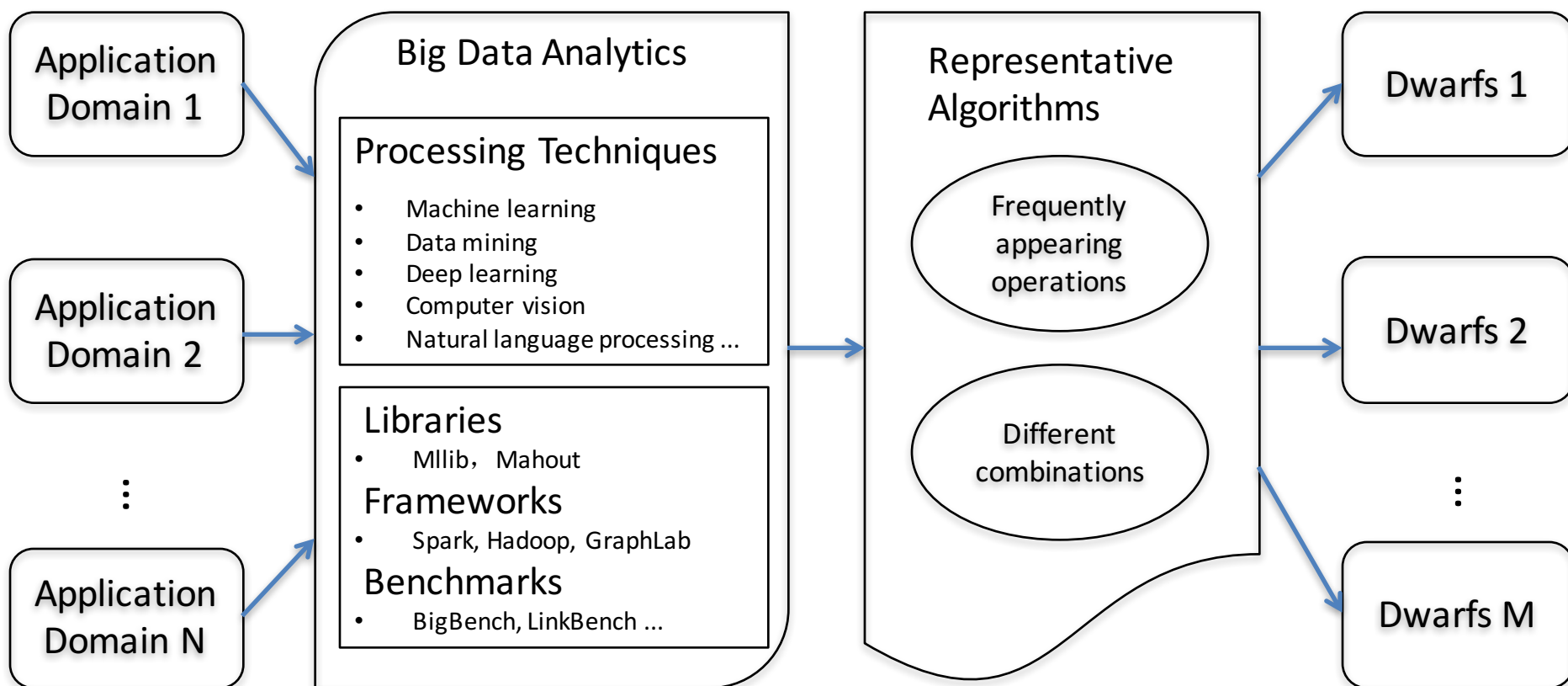
Frameworks

- Spark, Hadoop, Graphlab ...

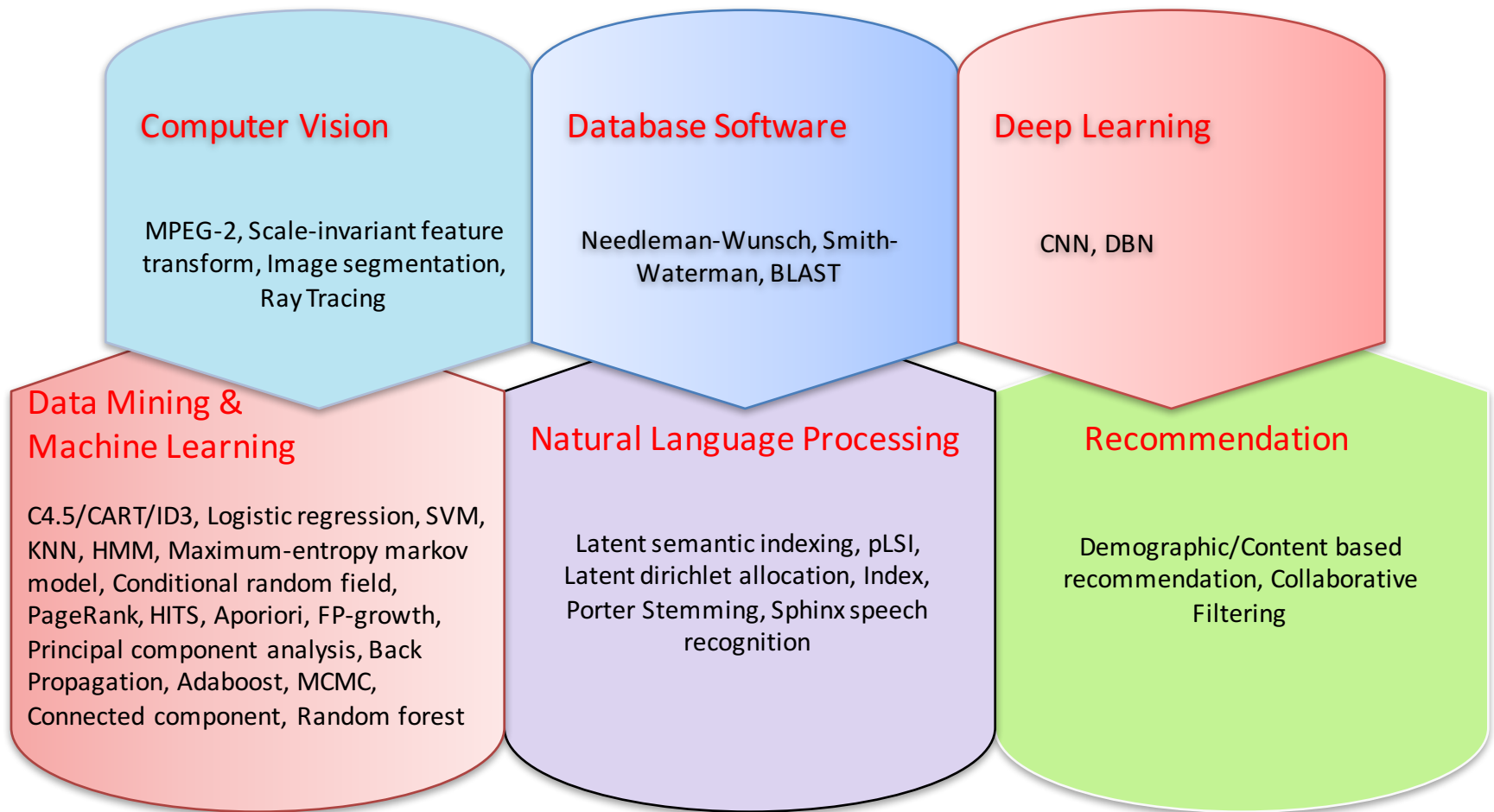
Benchmarks

- BigDataBench, Linkbench ...

Methodology of Dwarfs



Algorithms Chosen to Investigate



Frequently-appearing Operations

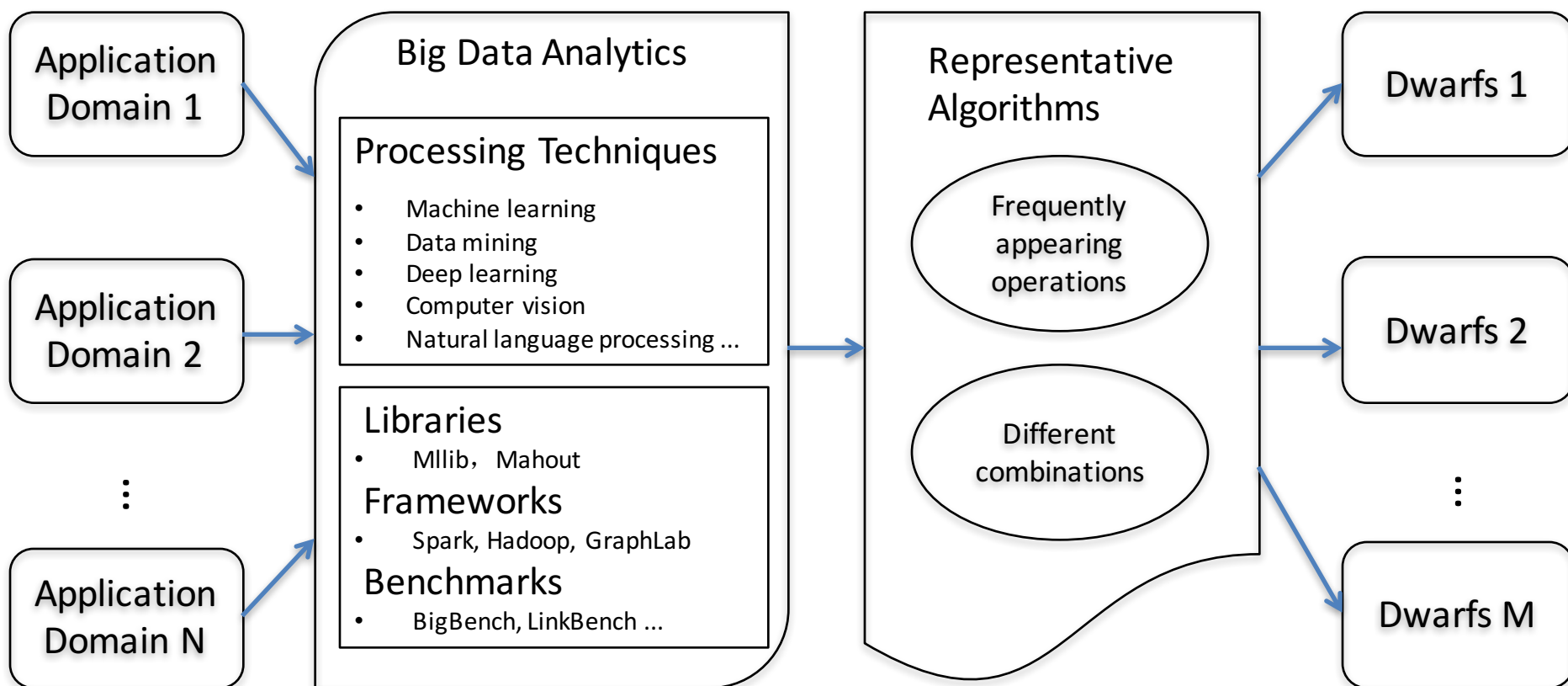
Statistic Operation

- Probability calculation
- LSI, pLSI, Latent dirichlet allocation ...

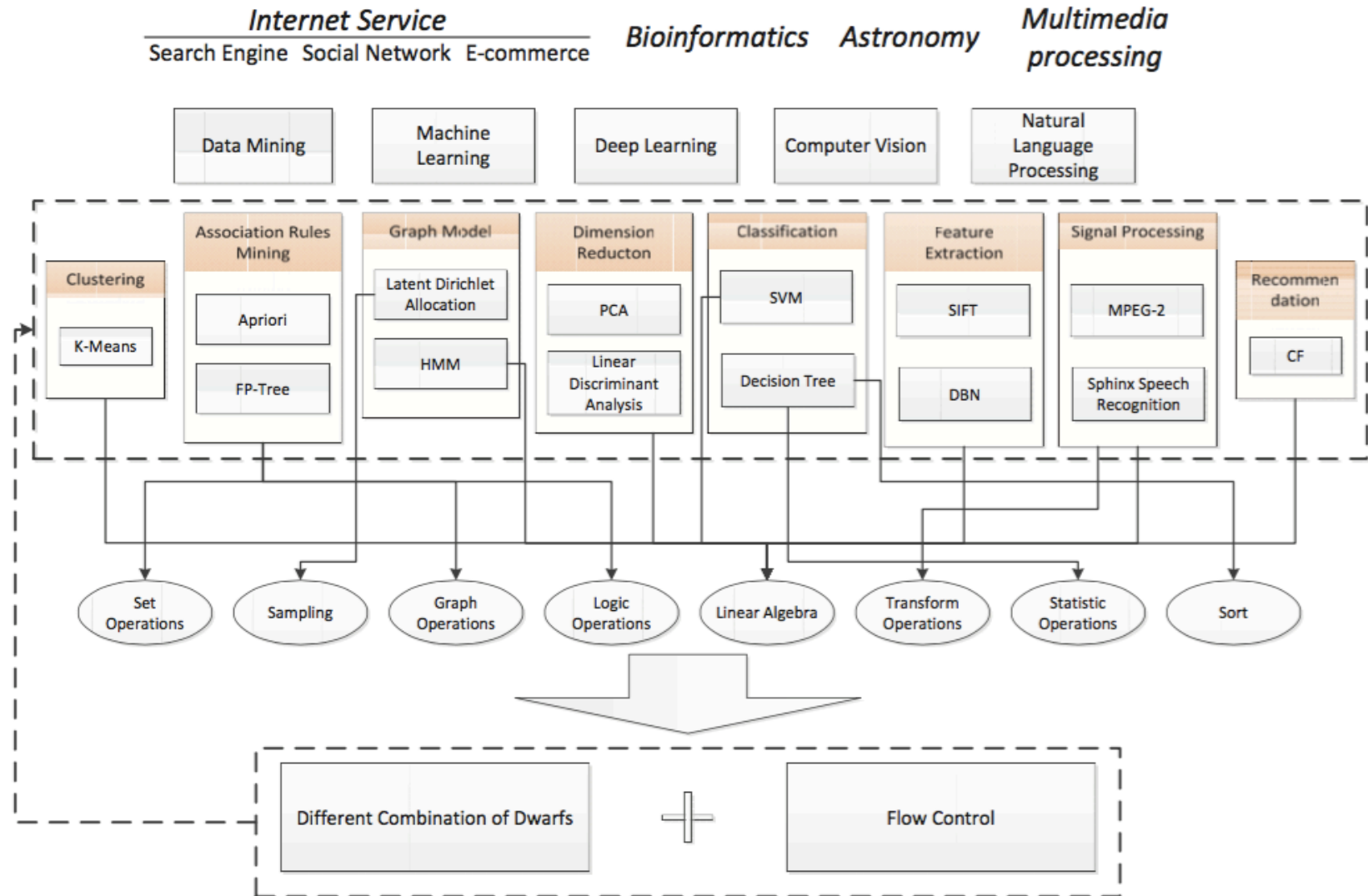
Sort

- Partial sort, quick sort, Top k sort...
- K-means, Decision tree ...

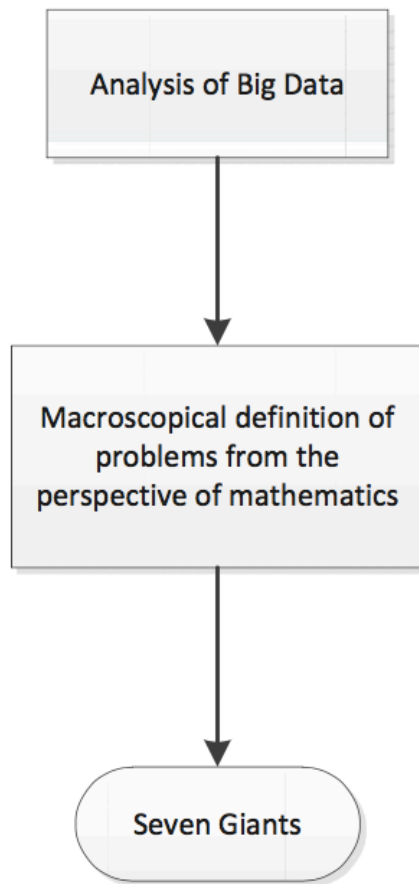
Methodology of Dwarfs



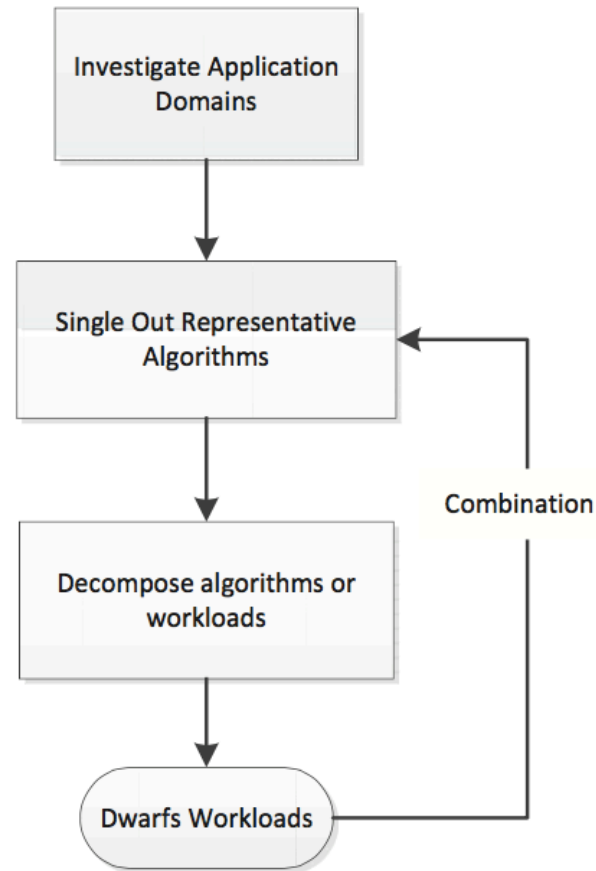
Methodology Summary



Comparison With NRC's Giants



(a) NRC Seven Giants



(b) Our Eight Dwarfs

- *N. Council. Frontiers in massive data analysis. The National Academies Press Washington, DC, 2013.*

Dwarfs in Big Data Offline Analytics

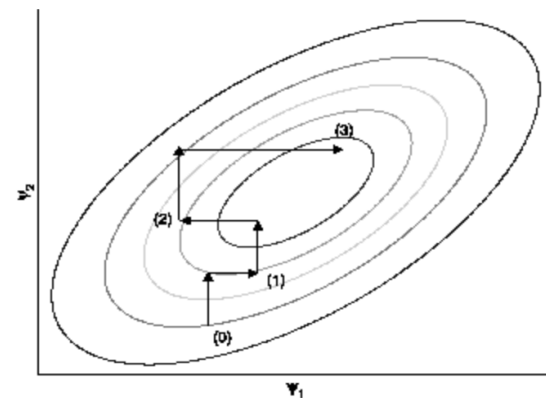
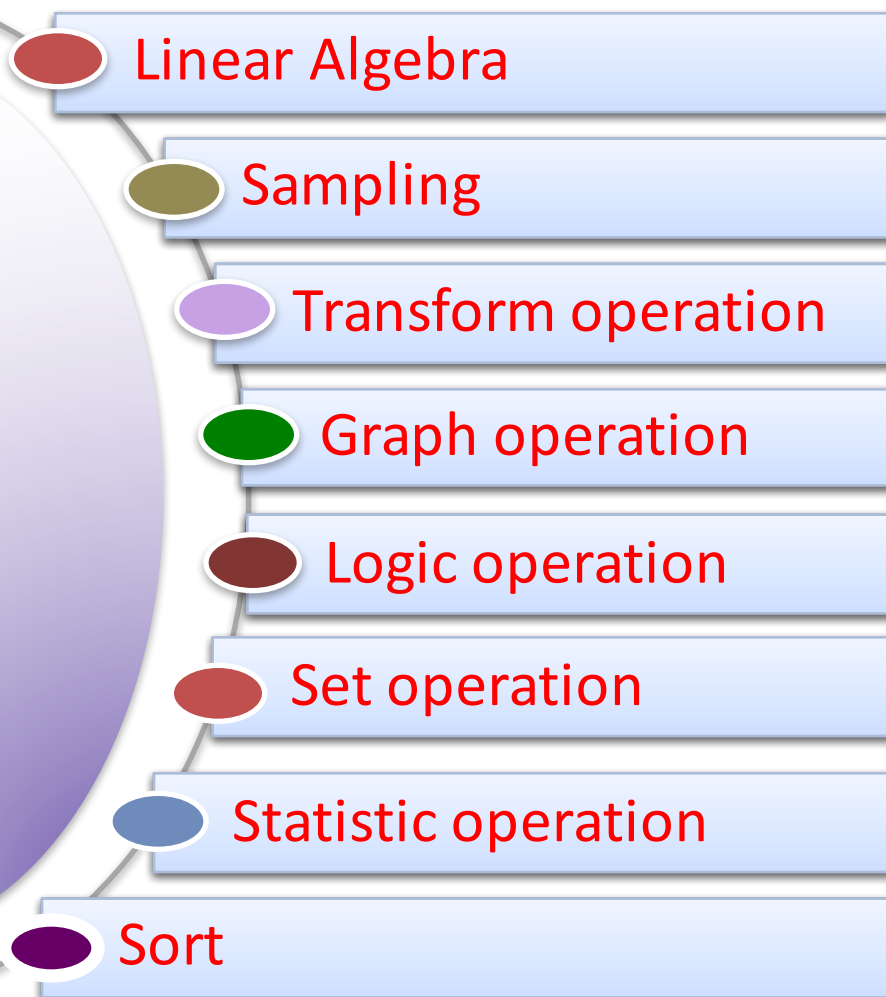
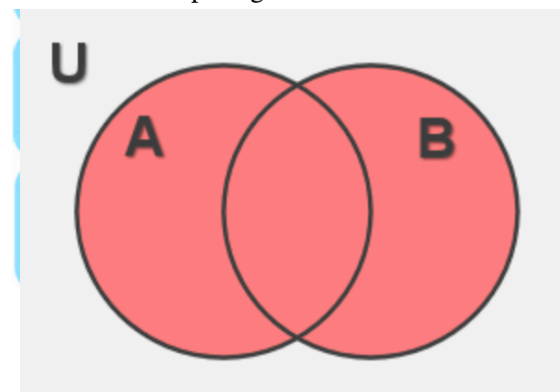


Figure 3.4: Gibbs sampling algorithm in two dimensions starting from an initial point and then completing three iterations

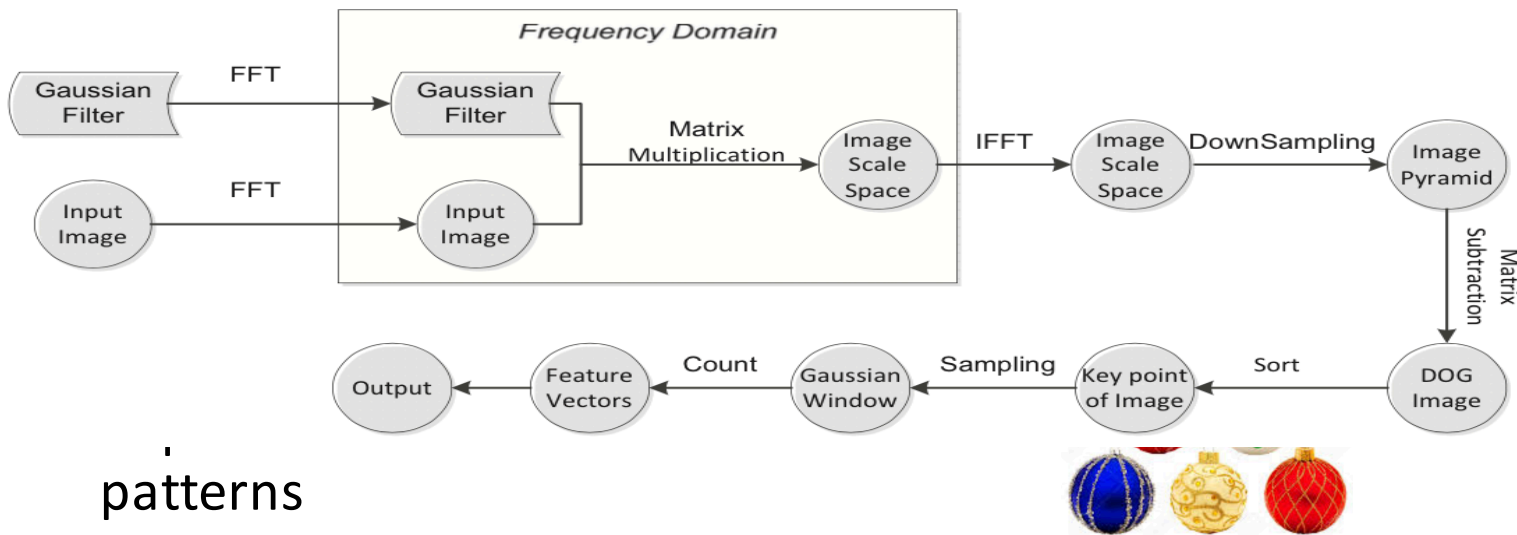


Properties

- Composability
 - One Algorithm can be composed of one or more dwarfs



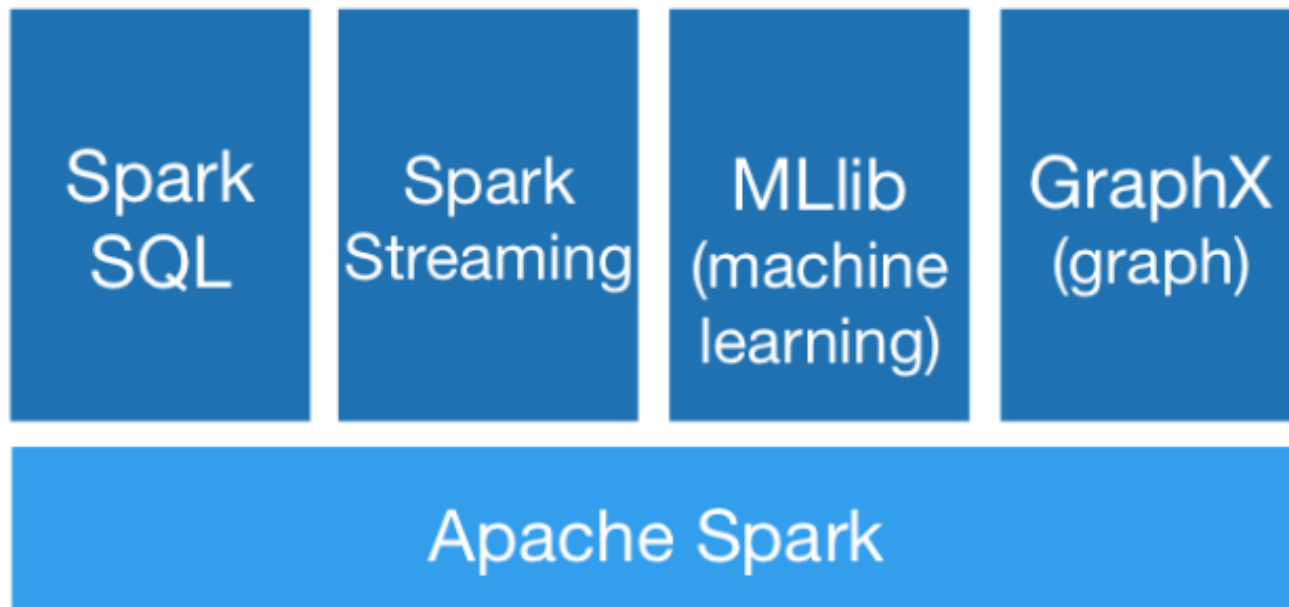
Feature extraction – SIFT Algorithm



patterns

Overlap? Yes

- View of applications
 - Same application can be implemented with different abstractions
- Provide a method for application optimization and finding bottleneck in different levels



Thank You!

Eight Basic Dwarfs Implementation

No.	Dwarfs	Implementation
Dwarf-1	Linear Algebra	Matrix Multiplication
Dwarf-2	Sampling	MCMC for GMM
Dwarf-3	Transform Operation	FFT
Dwarf-4	Graph Operation	BFS
Dwarf-5	Logic Operation	MD5
Dwarf-6	Set Operation	Union
Dwarf-7	Sort	QuickSort
Dwarf-8	Statistic Operation	Count