

FONDEMENT DES SYSTÈMES INFORMATIQUES

MapReduce

ARPEGE 2010

ANR-10-SEGI-001



Context and objectives

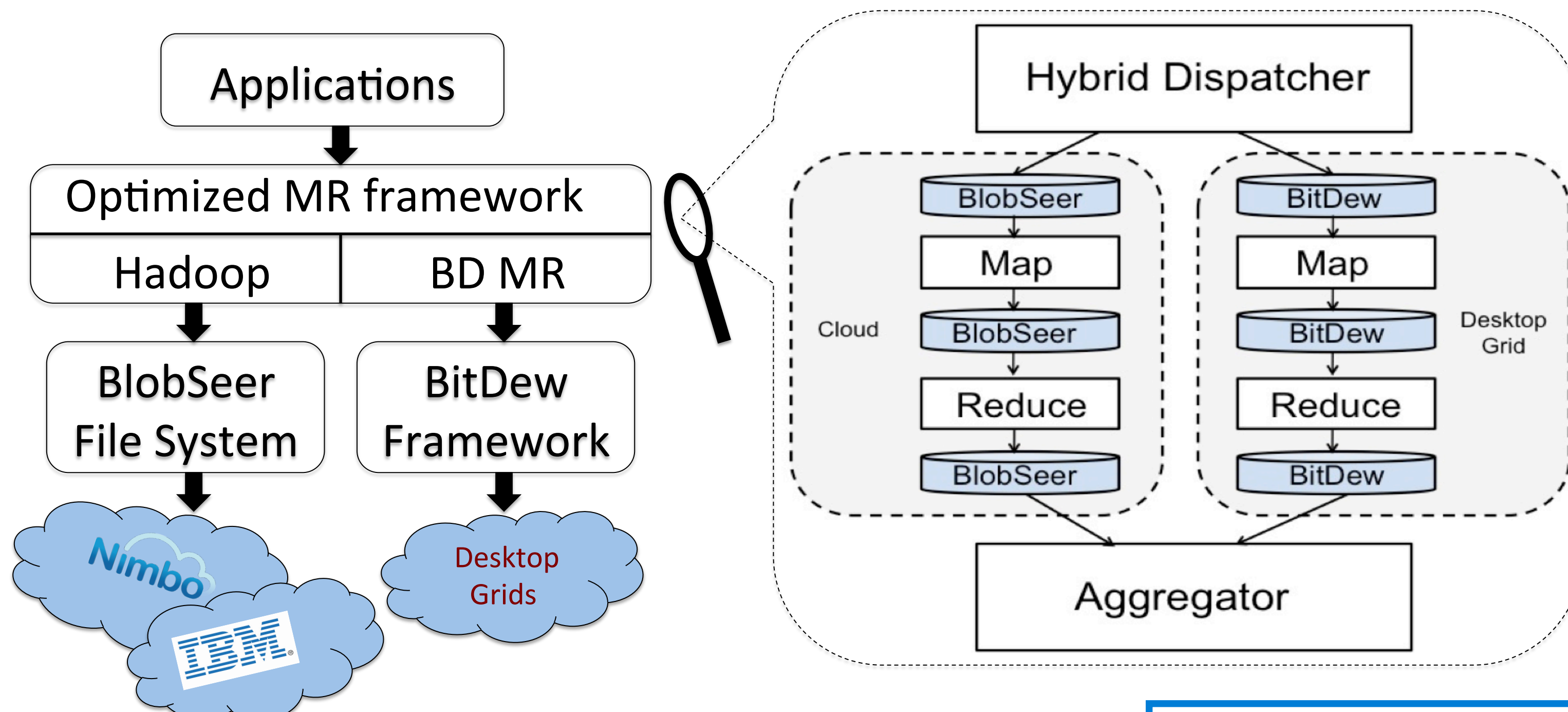
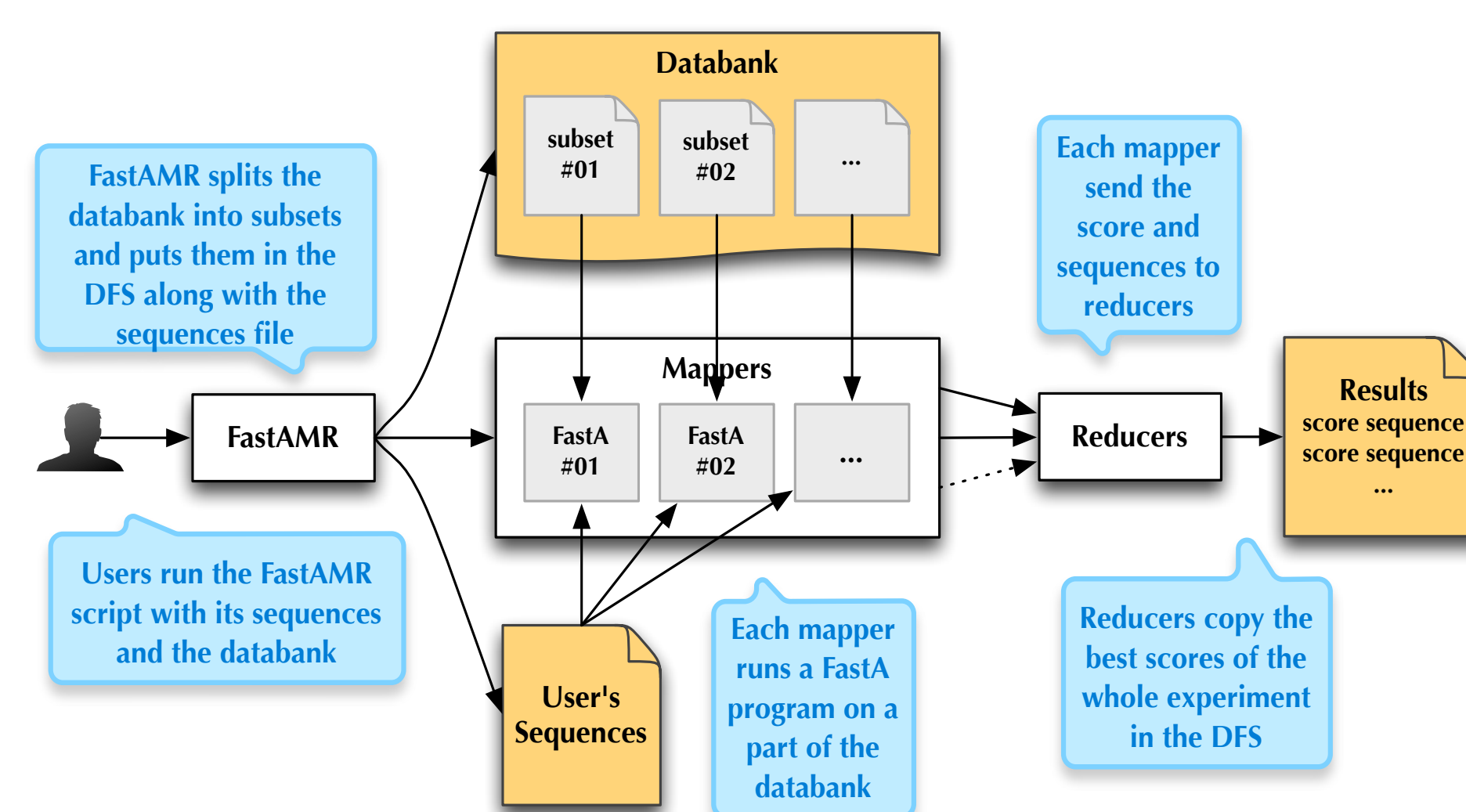
MapReduce is a programming model for data-intensive computing

Open issues and challenges:

- Low throughput for massively concurrent accesses
- Scheduling and fault tolerance still rudimentary
- Hybrid platforms (cloud federations, desktop grids) not explored yet

Goal: an optimized MapReduce framework for hybrid infrastructures

MapReduce of Biological Sequences

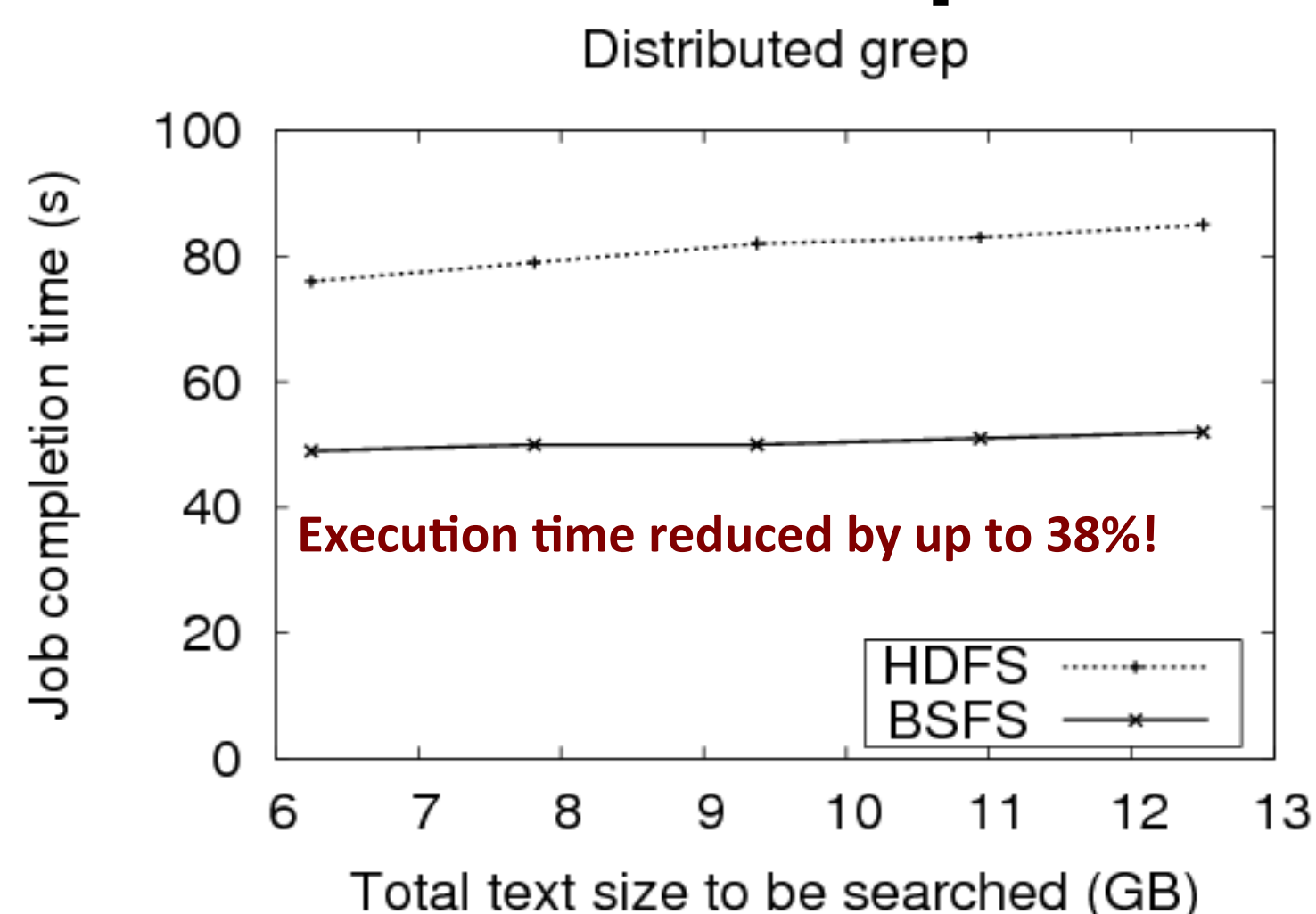


Methodology

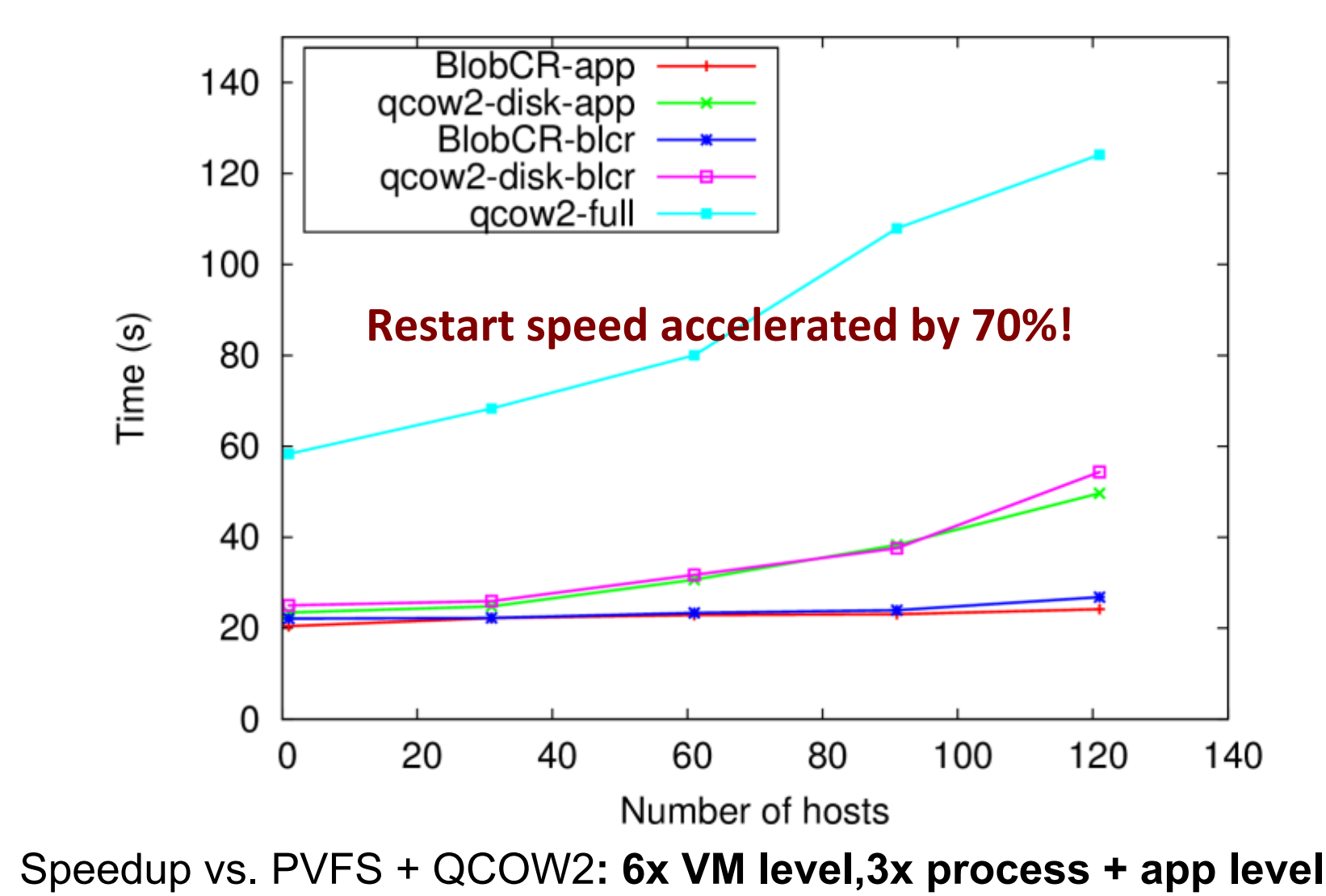
- Hybrid storage infrastructure
 - BlobSeer (BS): distributed storage management on Clouds
 - BitDew (BD): distributed storage on desktop grids
- High throughput concurrent data access
 - Distributed metadata, lock-free access to storage
- High level component model (HLCM)
 - Generic hierarchical connector-based component model

Highlights

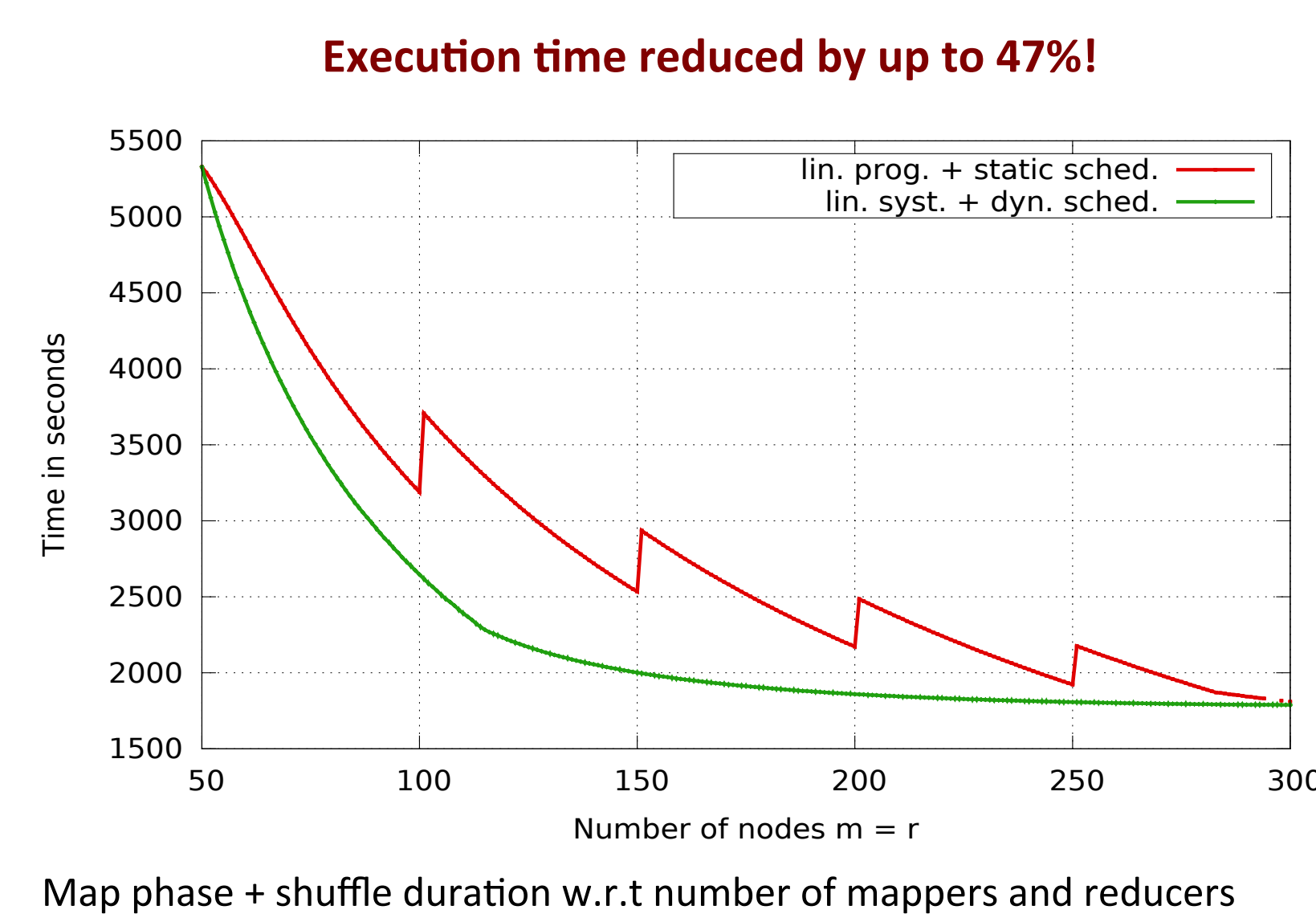
BlobSeer Does Better than Hadoop!



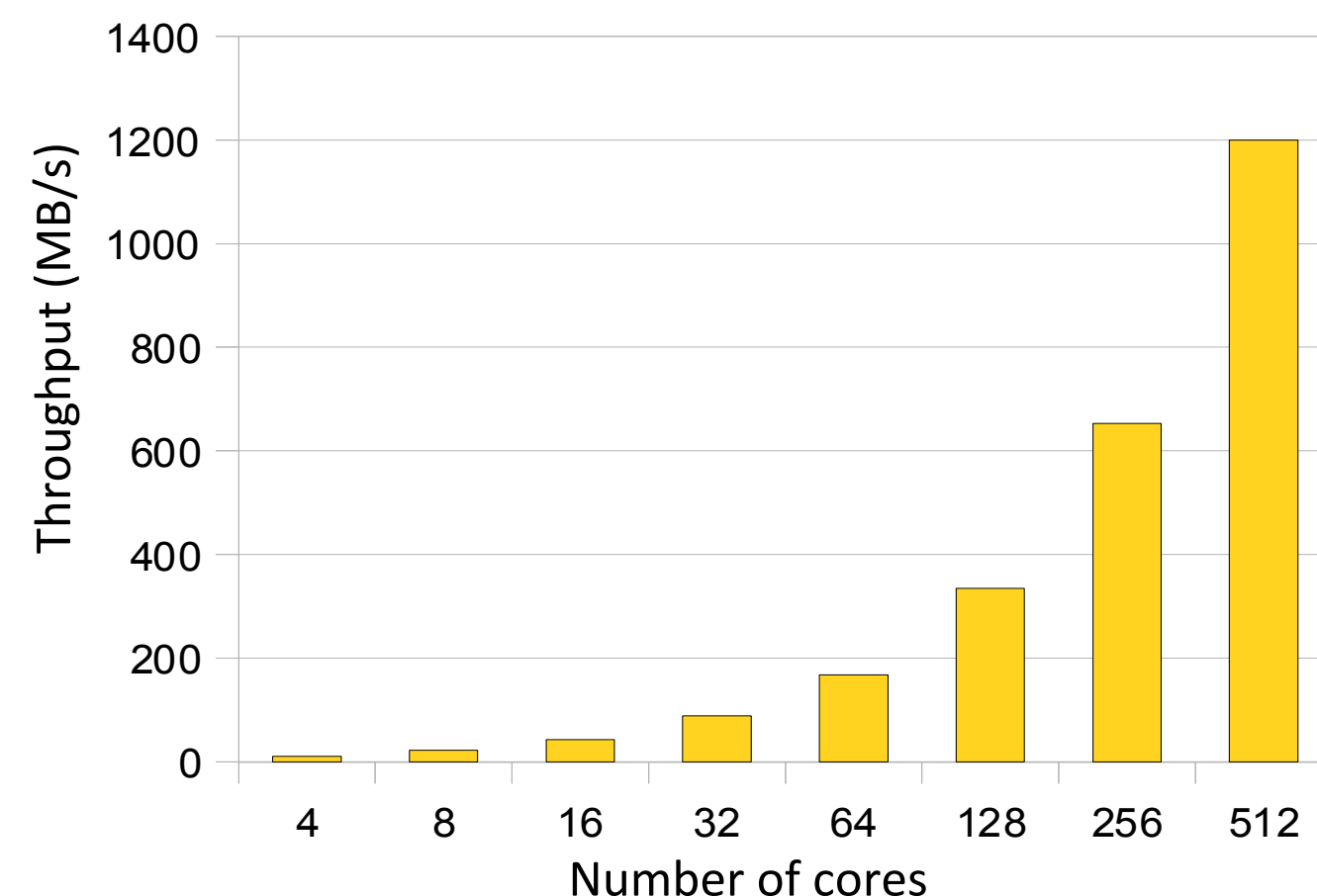
BlobCR: high performance resilience using virtual disk based checkpoint-restart



Scheduling



BitDew Scalable MapReduce Processing on Internet Desktop Grid

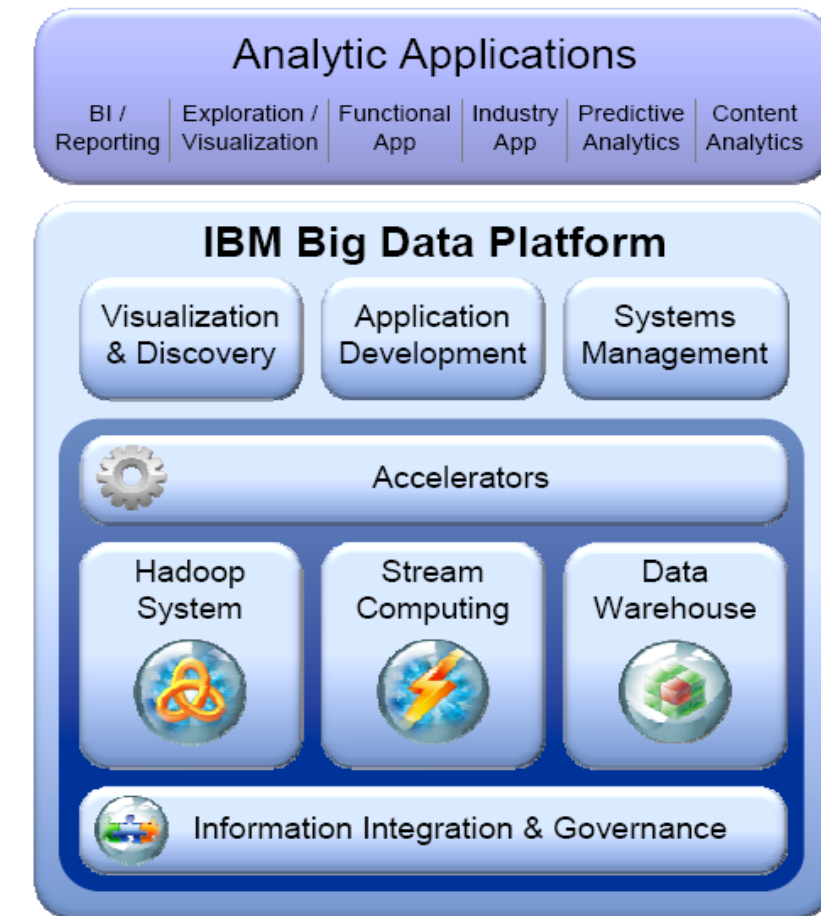


Applications



Bioinformatics

- SuMo – Protein Structure
 - Perform coarse-grain pairwise comparisons of protein structures
 - Compare a set of structures against a database
 - Commercial version: MEDIT
- FastA – Sequence Comparison
 - Compare sequences over a database
 - Compare a database over itself
 - Application to genome and proteome



COORDINATOR: Gabriel Antoniu, INRIA Rennes – Bretagne Atlantique

PARTNERS: INRIA (KerData, AVALON Project-Teams), CNRS IBCP, IBM France, Argonne National Lab, University of Illinois – Urbana Champaign, MEDIT

CONTACT :

gabriel.antoniu@inria.fr

More on MapReduce:
mapreduce.inria.fr



LES RENCONTRES DU NUMÉRIQUE

17 et 18 avril 2013