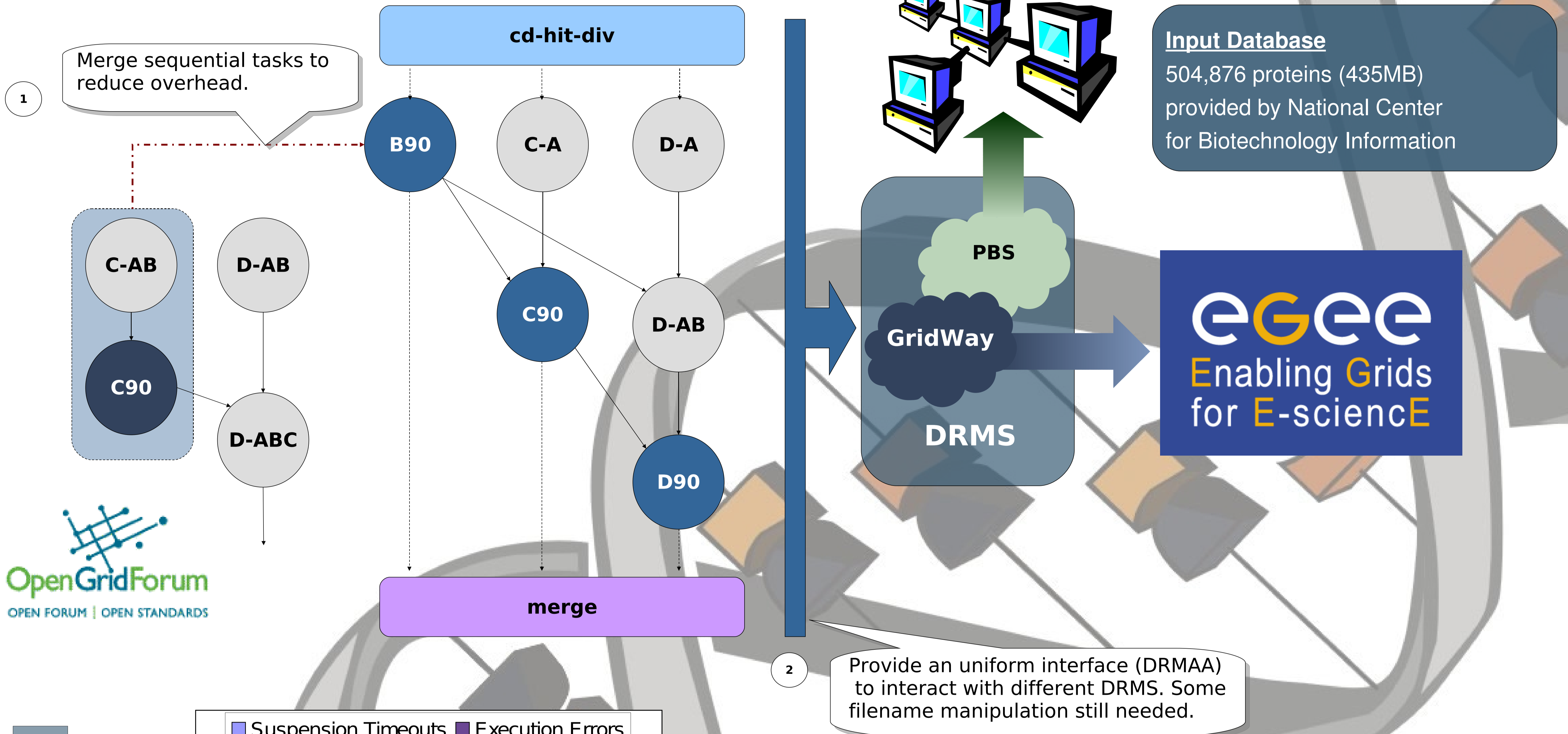


Protein Clustering

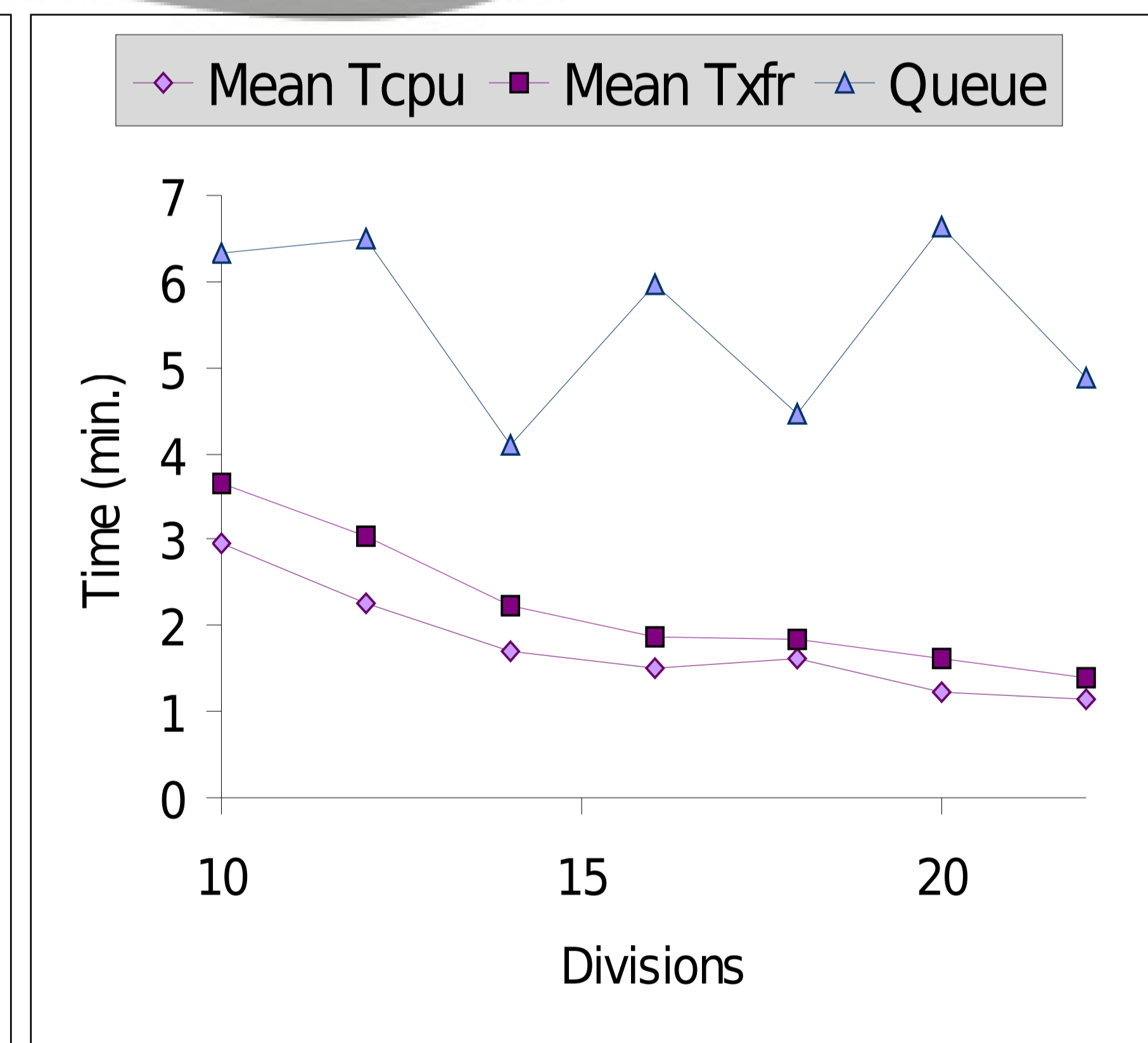
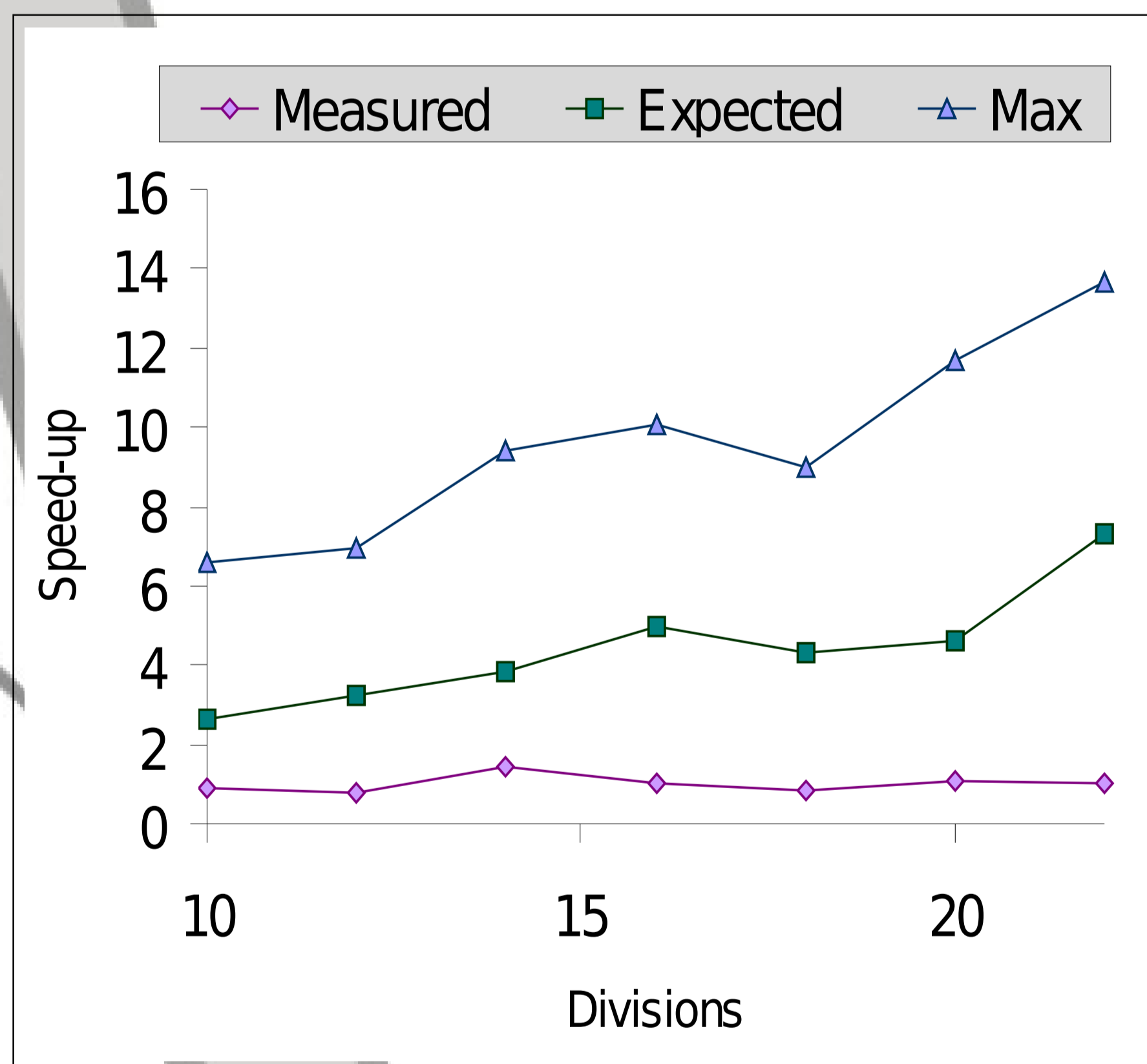
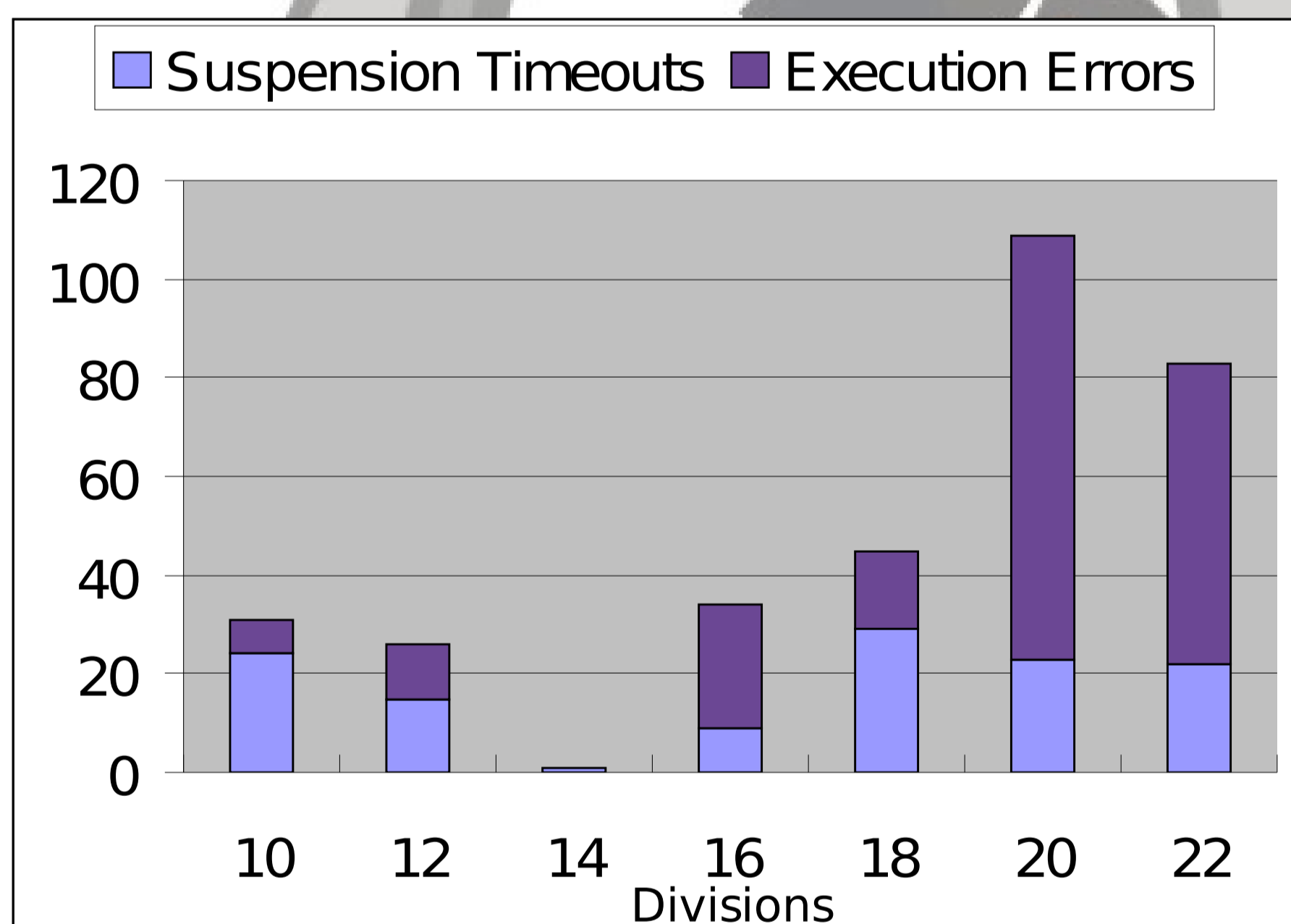
Process where protein databases are organized into groups or families in order to capture protein similarity. It can be applied in domain analysis, large protein database organization and search improvement.

CD-HIT

Toolkit for clustering large protein databases at high sequence identity threshold. Redundant sequences are removed and a database of only representatives is generated.



Results



GridWay (<http://www.gridway.org/>)

- Metascheduler standing on top of Globus services.
- Handles DAG based workflows.
- Allows advanced flow structures (loops, branches).
- Implements the Distributed Resource Management Application API (DRMAA) which is an OGF Standard.
- Considers static and dynamic resource information.
- Offers automatic staging mechanisms.
- Provides fault tolerance mechanisms (network outage, remote and local machine crash):
 - Tries task execution/file transfer on the same resource.
 - Submits failed task to an alternate resource.
 - Failed tasks are moved transparently to other resources.

Optimization

