# Protein Multiverse on University HPC Grid

## Azhar Ali Shah
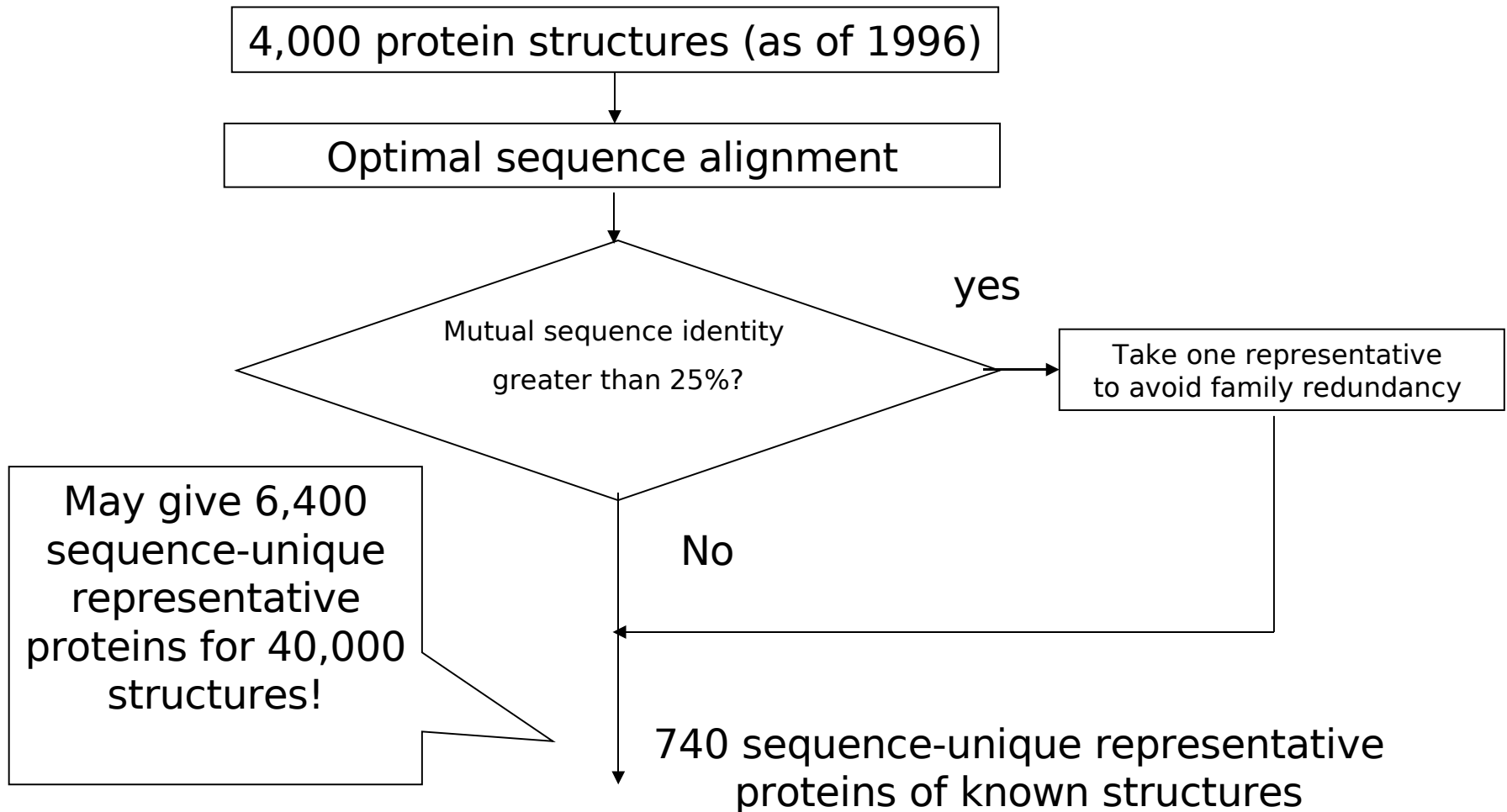
Protein Multiverse Meeting, Sep 26, 2007

# Outline

- **Related Work**
- **Complexity of the Problem**
- **Architectural Design**
- **Program Workflow Design (PWD)**
- **Infrastructure Details**
- **Discussion**

Azhar A Shah                    Protein Multiverse on University HPC Grid
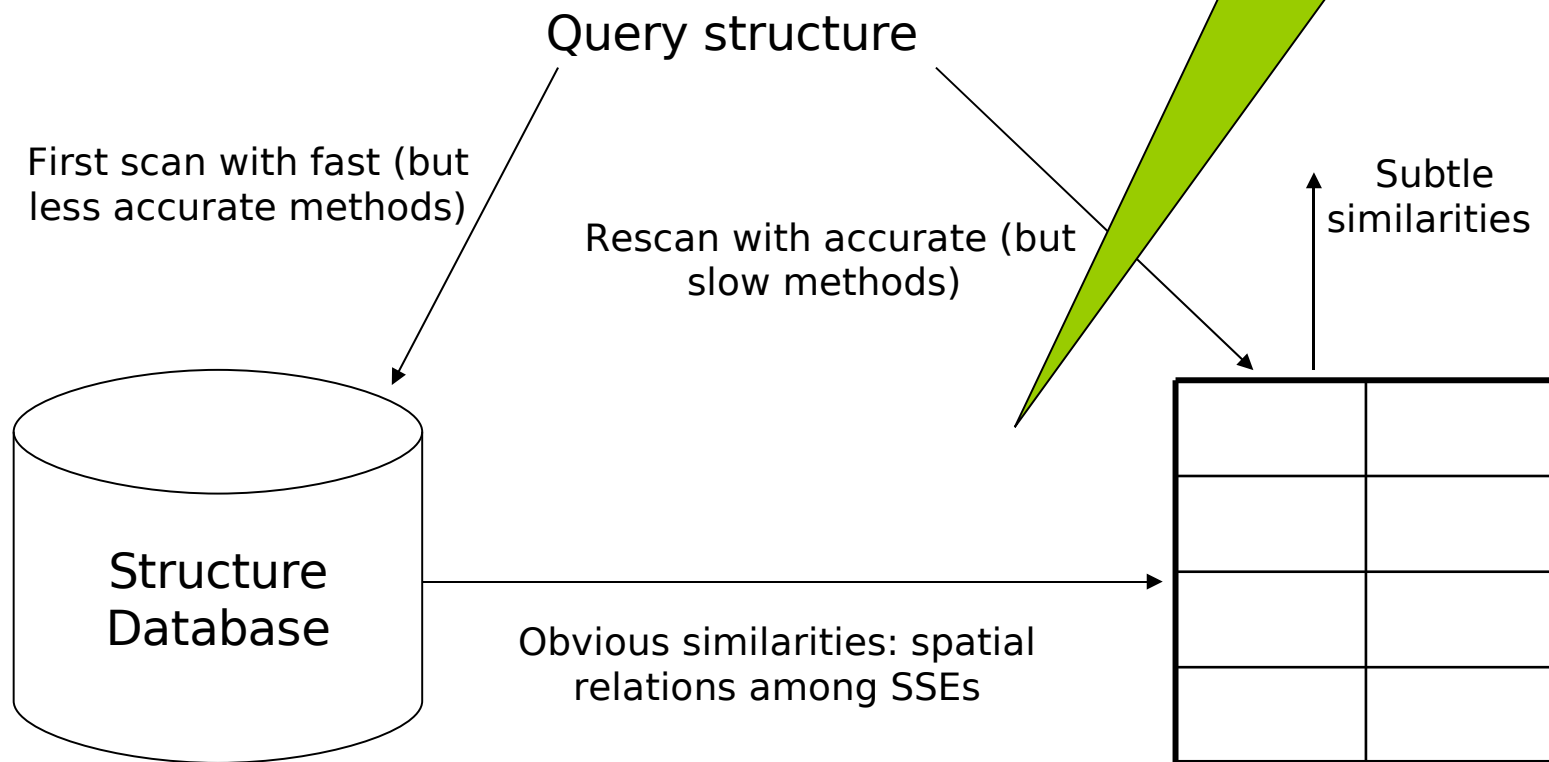
# Related Work1

□ Mapping the Protein Universe (Holm, Sander 1996)

  ■ **Motivations for all-on-all comparison:**
    □ Distribution of known structures in shape space
    □ Grand view of the architecture of all proteins
    □ A map of physical attractor regions in the abstract shape space of proteins
      ▪ Help to understand protein folding and evolution

# Database Preparation

4,000 protein structures (as of 1996)

Optimal sequence alignment

Mutual sequence identity greater than 25%?

yes

Take one representative to avoid family redundancy

No

May give 6,400 sequence-unique representative proteins for 40,000 structures!

740 sequence-unique representative proteins of known structures

Azhar A Shah                    Protein Multiverse on University HPC Grid

# Searching 3D Datab

**No parallelization**

Query structure

First scan with fast (but less accurate methods)

Rescan with accurate (but slow methods)

Subtle similarities

Structure Database

Obvious similarities: spatial relations among SSEs

One structure against several thousand structures takes 5 minutes on a normal workstation

Azhar A Shah                    Protein Multiverse on University HPC Grid

# Classification

On a high-dimensional fold space:

*Families:* Close range clusters
*Folds:* Intermediate range clusters
*Attractors:* Long range clusters

- Domains:
  - Structures having same recurring substructures are grouped into *Domains*
    - 1048 domains for 740 structures

- Fold classes:
  - Similar domains are grouped into *fold* classes
    - 287 folds for 1048 domains
  - Fold class is based on structural similarity and is analogous to **Family** which is based on sequence similarity.

- Attractors:
  - Five long regions in an abstract high-dimensional **fold space**

# Related Work2

- Global mapping of protein ***structure space*** and application in structure-based inference of protein function (Hou et al. 2005), PNAS.

  - Problem:
    - Simple structure comparison does not provide function inference for a protein with new fold

  - Solution:
    - A method based on map distance of protein structure space

Kim says. "This map provides us with a conceptual framework to organize ***all protein structures*** and functions and have that information readily available in one place", Berkley Lab Research News, Feb 2003.

# Database preparation

- PDB_SELECT 25 DATASET (Rel. Dec 2002)

  - A representative subset of the PDB containing *1,949* chains having <25% sequence identity

    - *51* chains further removed because of low resolution or length requirement of DaliLite

  - The resultant dataset consisted *1,898* protein chains

Azhar A Shah                    Protein Multiverse on University HPC Grid

# Mapping the protein structure space 1/3

IBM SP RS/6000

- The pairwise structural similarity of 1,898 chains was measured with DaliLite (**25,000 cpu hours)**

- The 1898x1898 similarity matrix $[s_{ij}]$ was converted to dissimilarity matrix $[d_{ij}]$ using:

$$dij = \begin{cases} s_{99.95} - s_{ij} \left( s_{99.95} > s_{ij}, i \neq j \right) \\ 0, \left( i = j \right) \\ s_{99.95, \left( otherwise \right)} \end{cases}$$

This matrix was used for structure space map (SSM)

Where $s_{99.95}$ is the 99.95th percentile of the distribution of all off-diagonal values of $s_{ij}$

Azhar A Shah                    Protein Multiverse on University HPC Grid

# Mapping the protein structure space 2/3

- Four scores:
  - Structure Space Map (SSM) distance score
  - DaliLite similarity score
  - DaliLite Z-score
  - BLAST-E values of pairwise sequence alignment

- ROC plots for evaluation
- Comparison of function inference among different scores based on GO function families
  - SSM outperforms other scores!

Azhar A Shah                    Protein Multiverse on University HPC Grid

# Mapping the protein structure space 3/3

▢ Based on present results it is predicted that the conceptual map of **all protein structures** would have same essential features.

Let us test this hypothesis!

Multi-method 3-D Map of Protein Structure Universe

Azhar A Shah                 Protein Multiverse on University HPC Grid

- Related work
- **Complexity of the Problem**
- Architectural Design
- Program Workflow Design (PWD)
- Infrastructure Details
- Discussion

# Complexity of the problem

▢ Job complexity

$$N_j = \frac{n(n-1)}{2} = \frac{41298 x 41297}{2} = 852741753$$
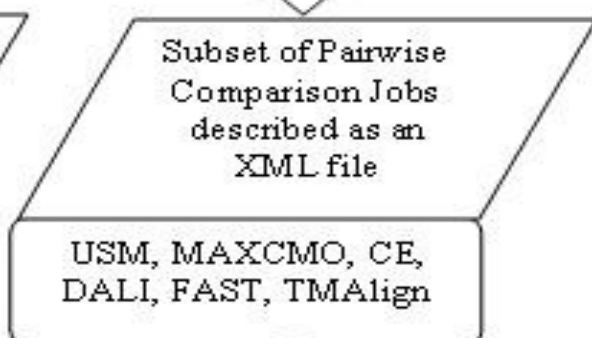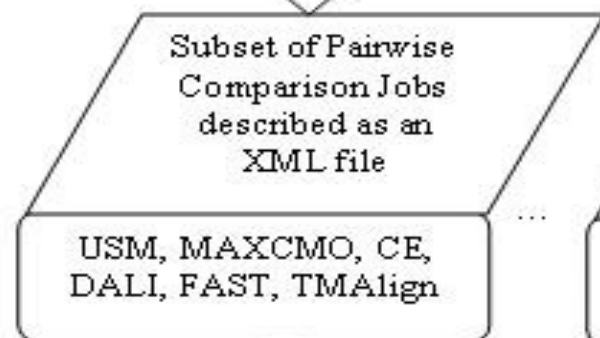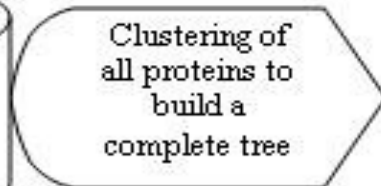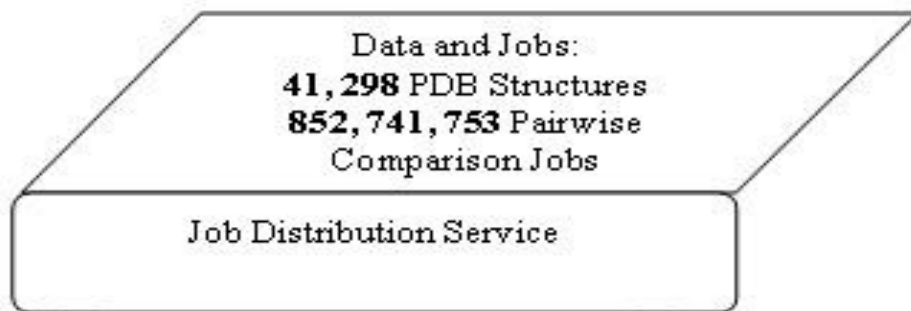
▢ Computational time
- 4088 hours => 170 days

▢ Storage complexity
- It takes 21 hours to download the PDB database with 41,298 structures which requires the space of 35 GB
- RAM would be the main obstacle for XML based input/output files

- Related work
- Complexity of the Problem
- **Architectural Design**
- Program Workflow Design (PWD)
- Infrastructure Details
- Discussion

Azhar A Shah                    Protein Multiverse on University HPC Grid

- Related work
- Complexity of the Problem
- Architectural Design
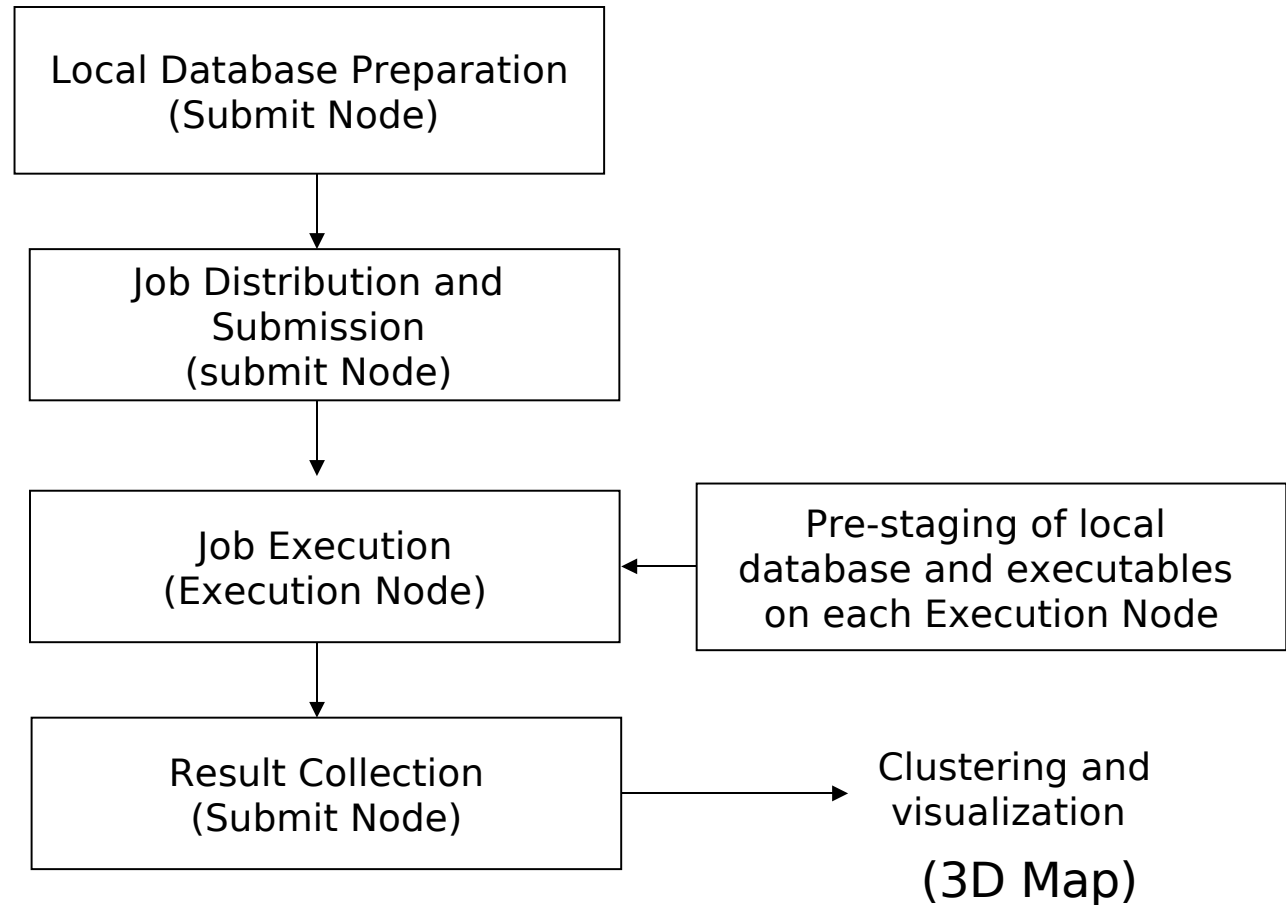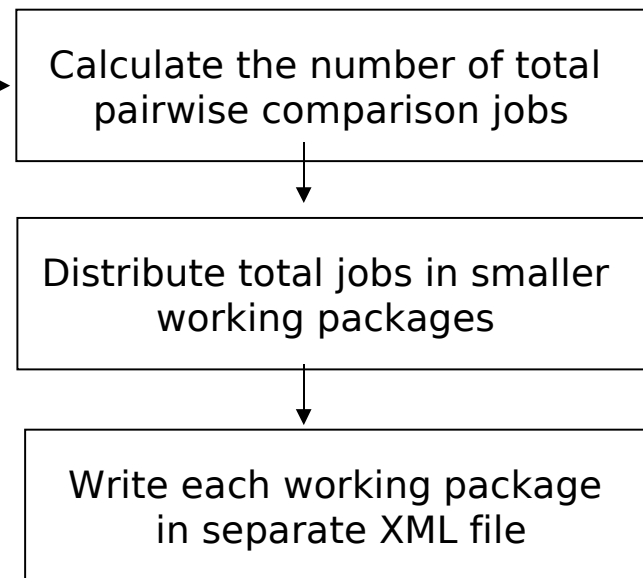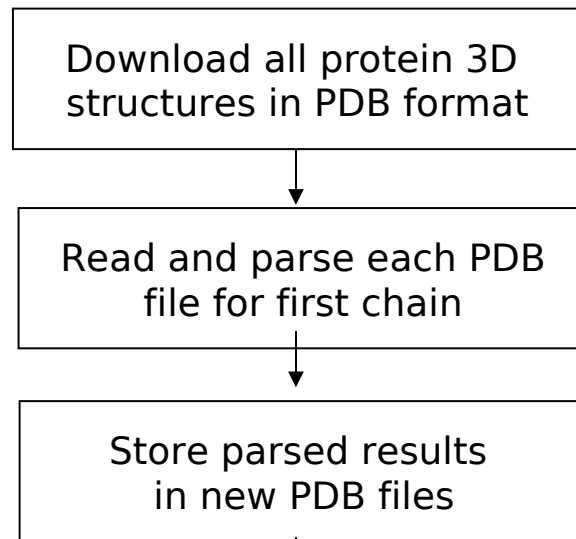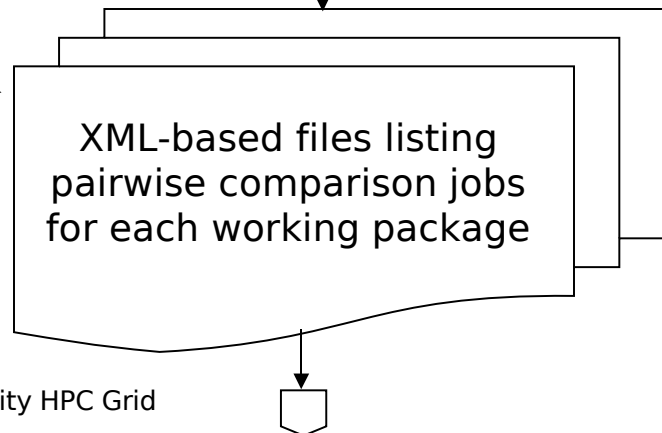- **Program Workflow Design (PWD)**
- Infrastructure Details
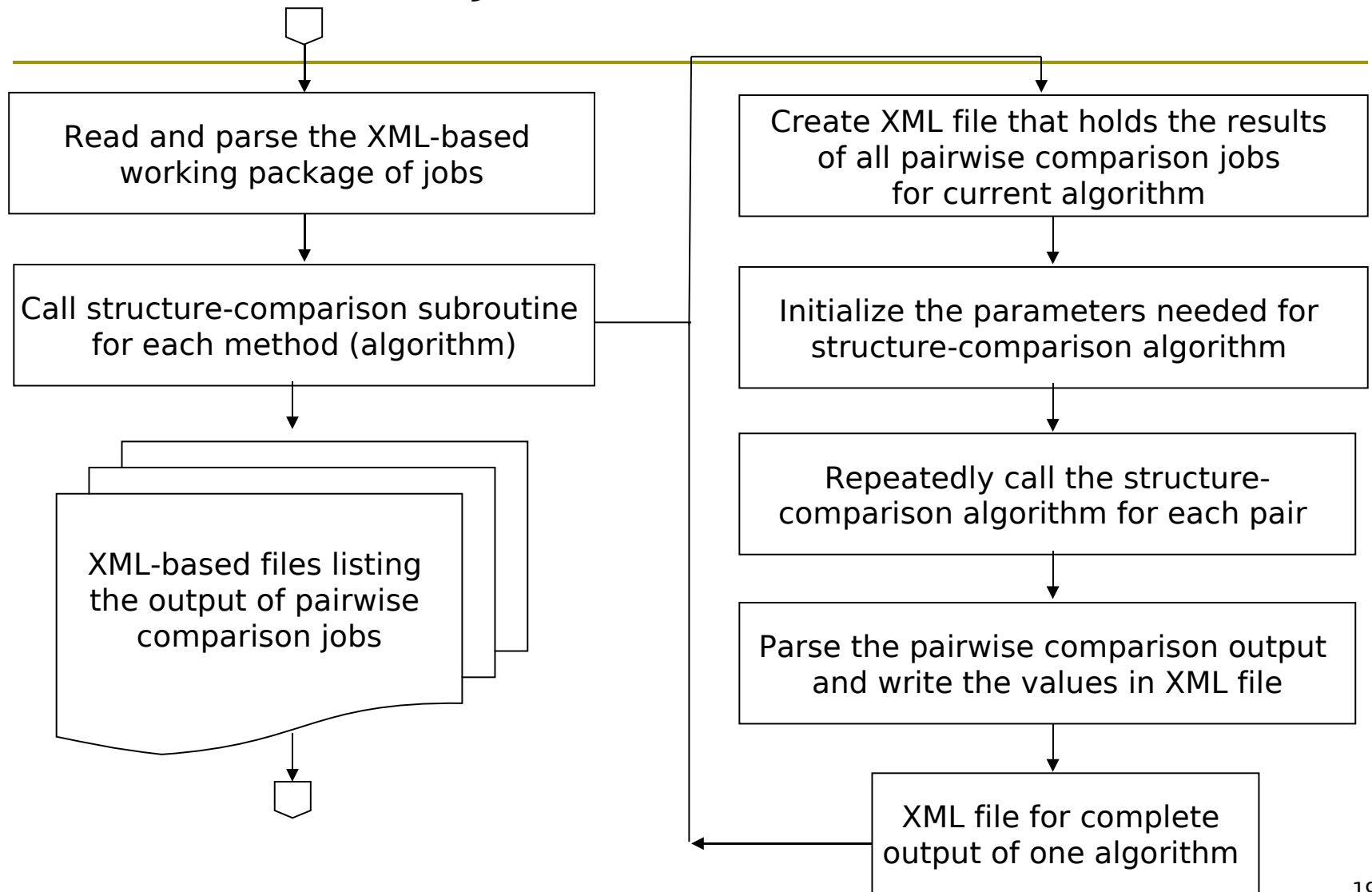- Discussion

Azhar A Shah                  Protein Multiverse on University HPC Grid

# PWD: Main Tasks



```
┌─────────────────────────────┐
│ Local Database Preparation  │
│      (Submit Node)          │
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│   Job Distribution and      │
│       Submission            │
│      (submit Node)          │
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐        ┌─────────────────────────────┐
│      Job Execution          │◀───────│    Pre-staging of local     │
│    (Execution Node)         │        │ database and executables    │
│                             │        │  on each Execution Node     │
└─────────────────────────────┘        └─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│      Result Collection      │───────▶  Clustering and
│       (Submit Node)         │           visualization
└─────────────────────────────┘             (3D Map)
```

Azhar A Shah                    Protein Multiverse on University HPC Grid

# Local Database Preparation

# Job Distribution

Download all protein 3D structures in PDB format

Read and parse each PDB file for first chain

Store parsed results in new PDB files

Calculate the number of total pairwise comparison jobs

Distribute total jobs in smaller working packages

Write each working package in separate XML file

```
<Work_Package id=1>
<Pair id=1  structure1=.... Structure2=.../>
<Pair id=2  structure1=....   structure2=.../>
        .
        .
        .
</Work_Package>
```

XML-based files listing pairwise comparison jobs for each working package

Azhar A Shah                    Protein Multiverse on University HPC Grid

# Job Execution



Read and parse the XML-based working package of jobs

Call structure-comparison subroutine for each method (algorithm)

XML-based files listing the output of pairwise comparison jobs

Create XML file that holds the results of all pairwise comparison jobs for current algorithm

Initialize the parameters needed for structure-comparison algorithm

Repeatedly call the structure-comparison algorithm for each pair

Parse the pairwise comparison output and write the values in XML file

XML file for complete output of one algorithm

Protein Multiverse on University HPC Grid

# Typical output XML file

```xml
<?xml version="1.0" encoding="U..."

<Method Name="CE">
  <Pair No="0" Structure1="1NTR...
    <Measures Align="124" RMSD...
  </Pair>
  <Pair No="1" Structure1="1NTR_-1.PDB" Structure2="1HTIA-1.PDB">
    <Measures Align="90" RMSD="3.93" ZScore="2.8" Se-ID="10.0" />
  </Pair>
```

```xml
<Work_Package id=1>
<Method name="CE">
  <Pair id=1  structure1=.... Structure2=...>
    <Measures Align=... ....... ...... />
  </Pair>

  <Pair id=2  structure1=....    structure2=.../>
    <Measures Align=... ....... ...... />
  </Pair>  .
</Work_Package>
```

Azhar A Shah                Protein Multiverse on University HPC Grid

# Result Collection



Read all the XML files containing the
results of each working package

Parse the method (algorithm) specific pairwise
comparison values and write them in the database

Maintain the
database
For all pairwise
comparison
results

Clustering and
visualization

Azhar A Shah                    Protein Multiverse on University HPC Grid

# Results Database Schema

| CE | | | |
|---|---|---|---|
| Pair_Lablel | RMSD | Z-Score | ... |
| Str1:Str2 | 123 | 123 | ... |
| ... | ... | ... | ... |

Azhar A Shah                     Protein Multiverse on University HPC Grid

- **Related work**
- **Complexity of the Problem**
- **Architectural Design**
- **Program Workflow Design (PWD)**
- **Infrastructure Details**
- **Discussion**

Azhar A Shah                    Protein Multiverse on University HPC Grid

# User Interface: Example

Azhar A Shah                    Protein Multiverse on University HPC Grid        www.google.com

# ProCKSI User Interface: Grid-based Portal Environment

Classical HTTP Interface
(for users working with
<100 pdb structures)

New FTP Interface
(for users working with
>100 pdb structures)

**Portlets:** Pluggable UI
components based on
web service architecture

Protein Multiverse
Interface

Workflow Design
Interface

Help and Support

Documentation and
Literature

Authenticated access for
collaborators/developers

Structure and Function
Prediction

???

Azhar A Shah                    Protein Multiverse on University HPC Grid

**Example: PROGRESS Portal Access:**
(Bogdanski Maciej et al. 2004)
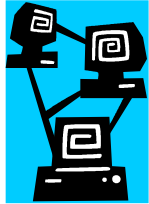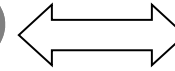
Azhar A Shah

Other future collaborators

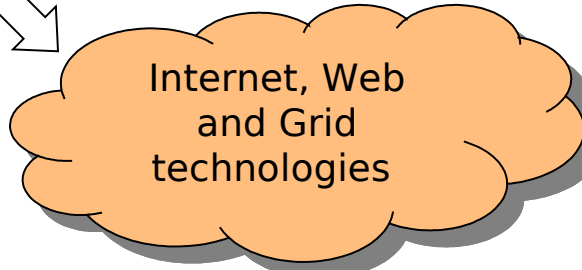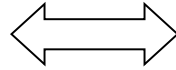University of Sindh, Pakistan

Azhar with so many dedicated PhD students to work on future trends for ProCKSI

Internet, Web and Grid technologies

University of Nottingham, UK

Protein multiverse experiments and overall grid architecture (multiverse portlet service)

Poznan Super Computing and Networking Centre, Poland

**Jacek** with one dedicated PhD student to develop Grid-based Portal Interface for ProCKSI

University of Calabria, Italy

**Gianluigi** with one dedicated PhD student to develop Grid-based workflow portlet service for ProCKSI

27

Azhar A Shah

Protein Multiverse on University HPC Grid

# University of Nottingham: Triton

Azhar A Shah                    Protein Multiverse on University HPC Grid

# University of Nottingham: Jupiter

Azhar A Shah                    Protein Multiverse on University HPC Grid

- **Related work**
- **Complexity of the Problem**
- **Architectural Design**
- **Program Workflow Design (PWD)**
- **Infrastructure Details**
- **Discussion**

Azhar A Shah                    Protein Multiverse on University HPC Grid

# Discussion

- ☐ Is it OK?

Azhar A Shah                    Protein Multiverse on University HPC Grid