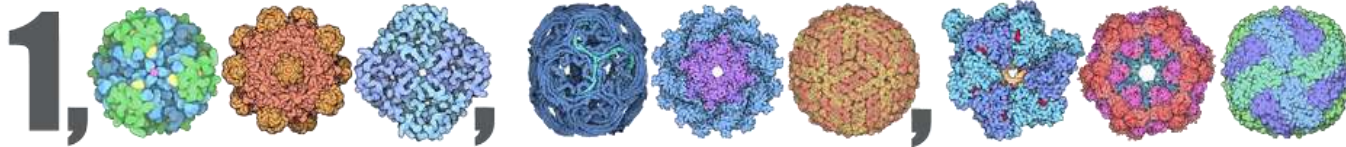


# Compressive Structural Bioinformatics: Large-scale analysis and visualization of the Protein Data Bank archive

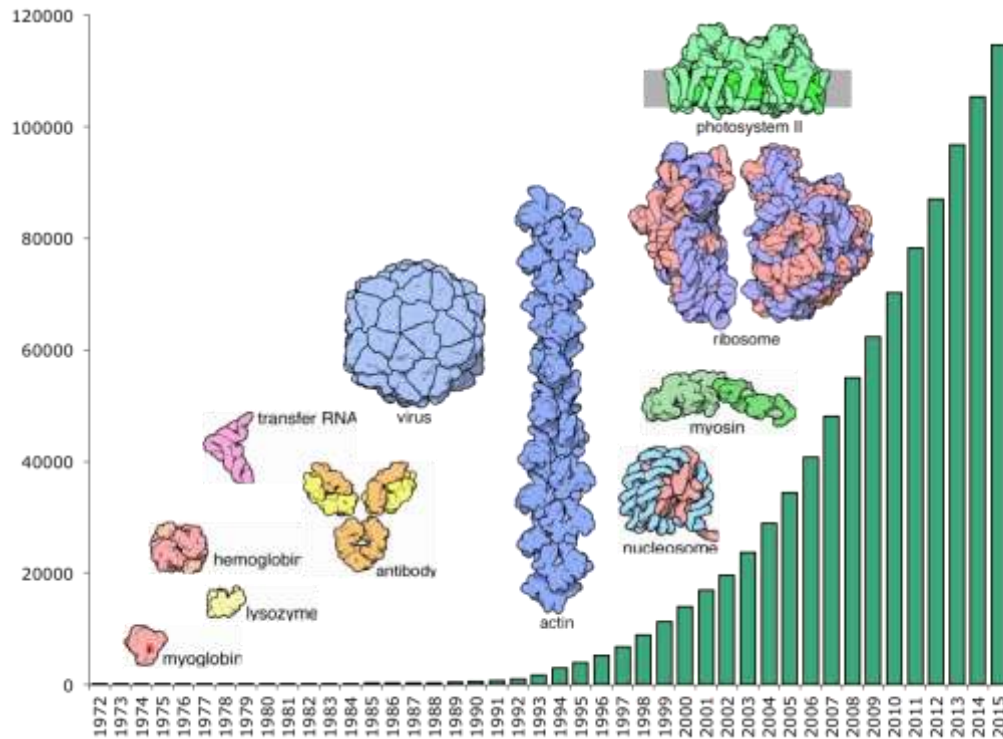
Peter W. Rose, Anthony R. Bradley,  
Alexander S. Rose, Yana Valasatava,  
Jose M. Duarte, Andreas Prlić

*Structural Bioinformatics Laboratory  
San Diego Supercomputer Center  
UC San Diego*

# PDB – A Billion Atom Archive

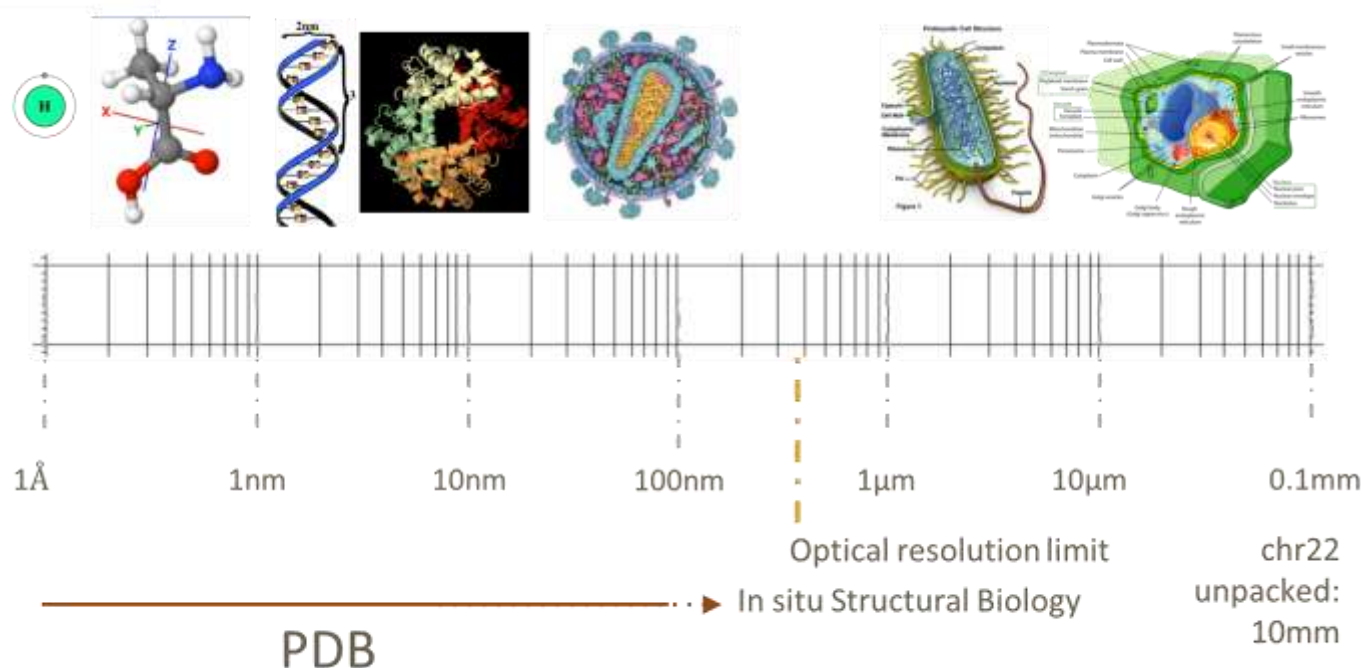


> 1 billion atoms in the asymmetric units

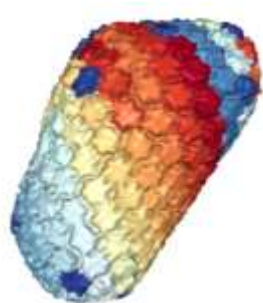


120,000  
structures  
in June 2016

# Growing Structure Size and Complexity

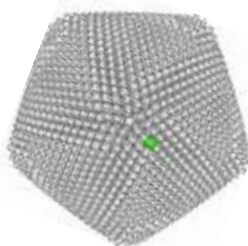


Largest asymmetric structure in PDB



HIV-1 capsid: PDB ID 3J3Q  
~2.4M unique atoms

Largest symmetric structure in PDB

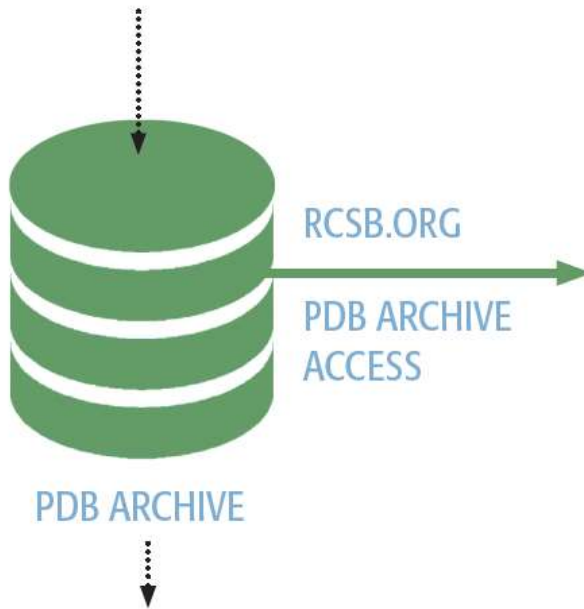


Faustovirus major capsid: PDB ID 5J7V  
~40M overall atoms

# Growing User Base

## ACCESSING PDB AND RCSB PDB

In 2015, 9329 entries were released into the PDB archive.



Total PDB archive traffic from all wwPDB partners totaled 534,339,871 downloads



Each month in 2015, **rccb.org** was visited 741,000 times on average by 315,000 unique visitors

A total of 35,260 GB of data were accessed



# → Scalability Issues

- **Interactive visualization**

- slow network transfer
- slow parsing
- slow rendering



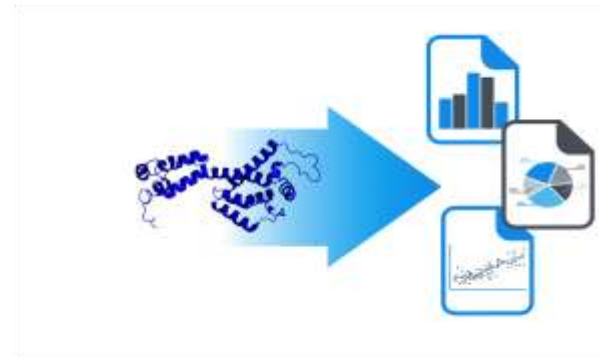
- **Mobile visualization**

- limited bandwidth
- limited memory



- **Large-scale structural analysis**

- slow repeated I/O
- slow repeated parsing

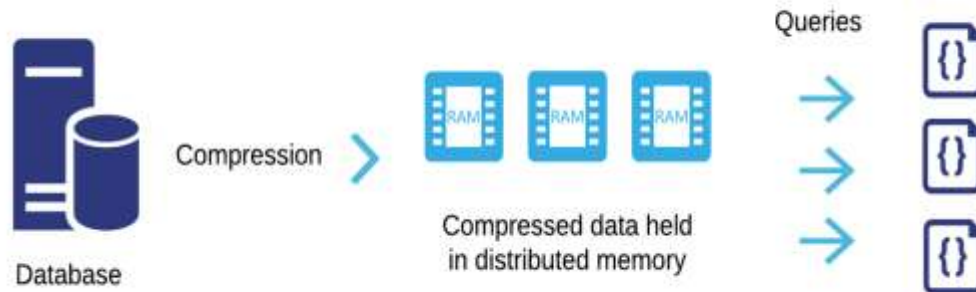


# Compressive Structural Bioinformatics

Efficiently store, transmit, and visualize 3D structures of biological macromolecules

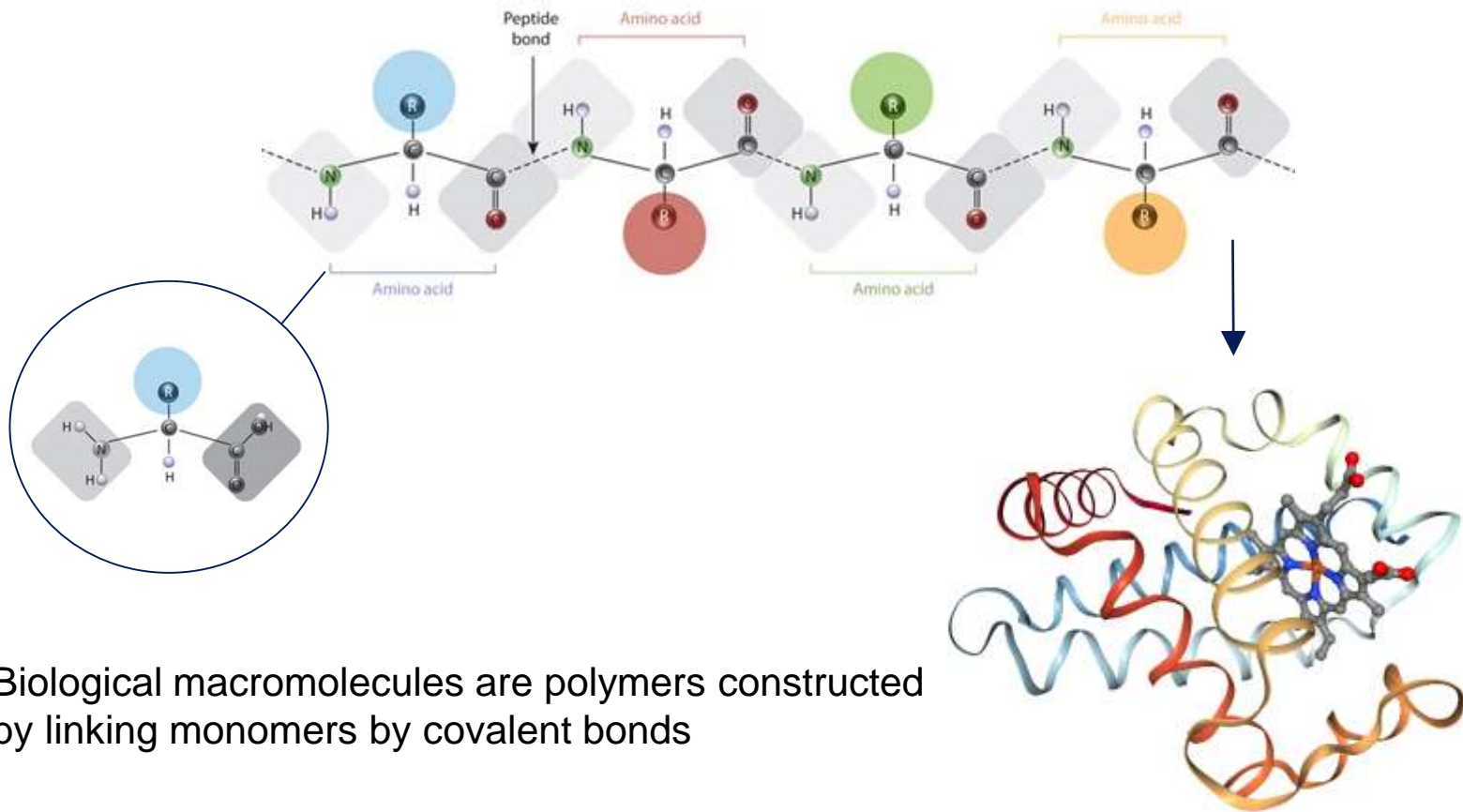


Perform large-scale structural calculations such as geometric queries or structural comparisons over the entire PDB archive held in memory



# Macromolecular 3D Structure

Biological macromolecules: proteins, nucleic acids



Biological macromolecules are polymers constructed by linking monomers by covalent bonds

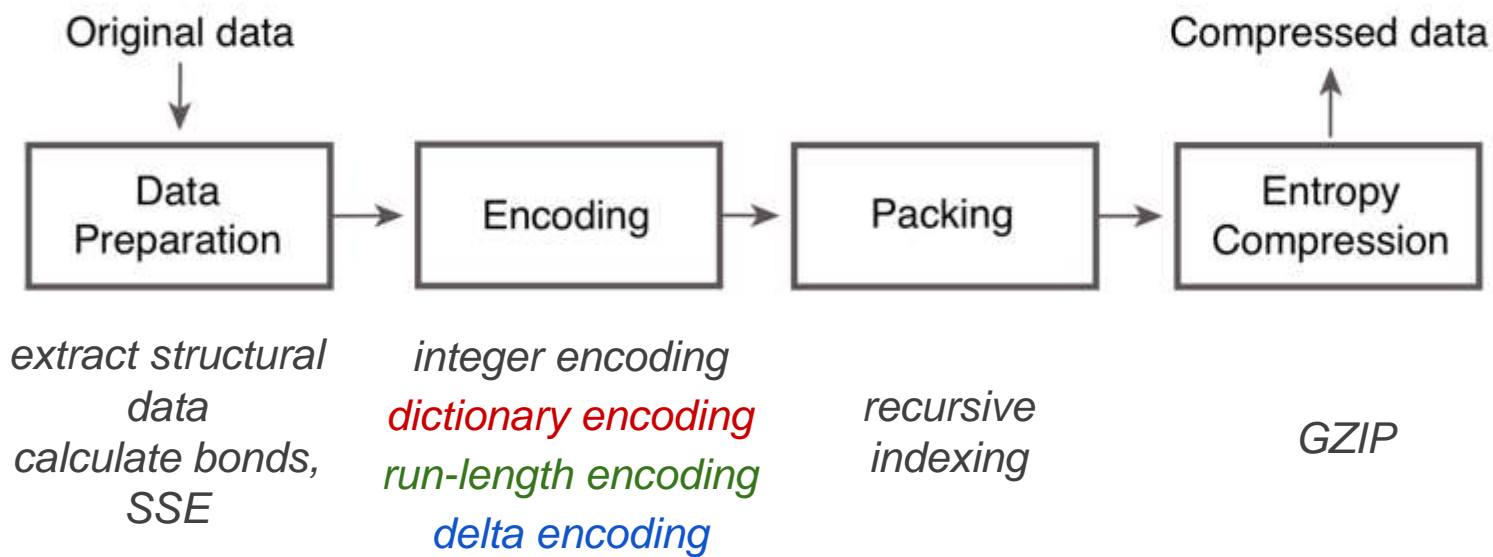






- **MacroMolecular Transmission Format ([mmtf.rcsb.org](http://mmtf.rcsb.org))**
  - Compact
    - fast network transfer, less I/O
  - Fast to parse
    - binary, no string parsing
  - Contains information for structural analysis and visualization
    - covalent bonds and bond orders
    - consistently calculated secondary structure

# MMTF Compression Pipeline



*Binary, extensible container format of MMTF*

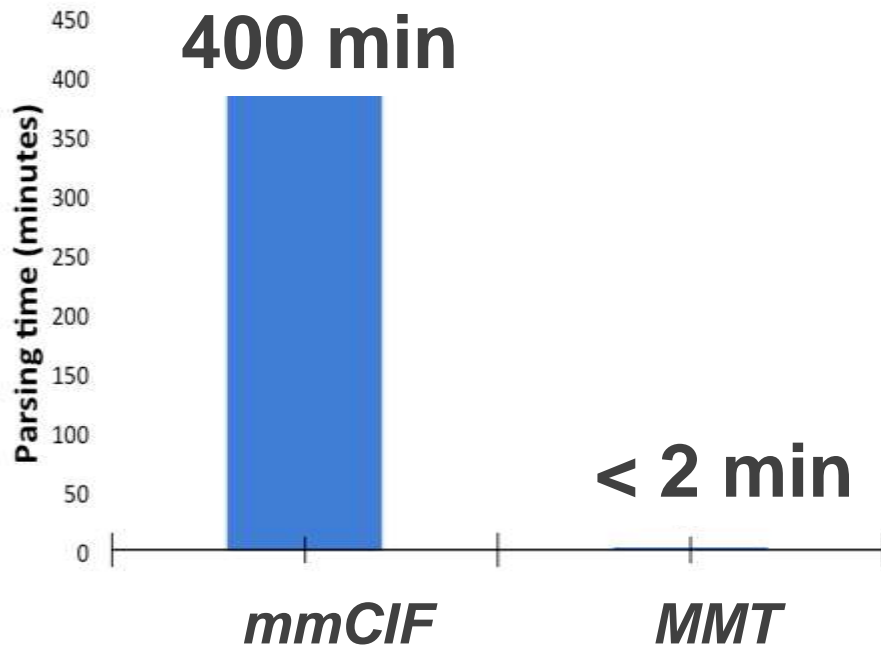
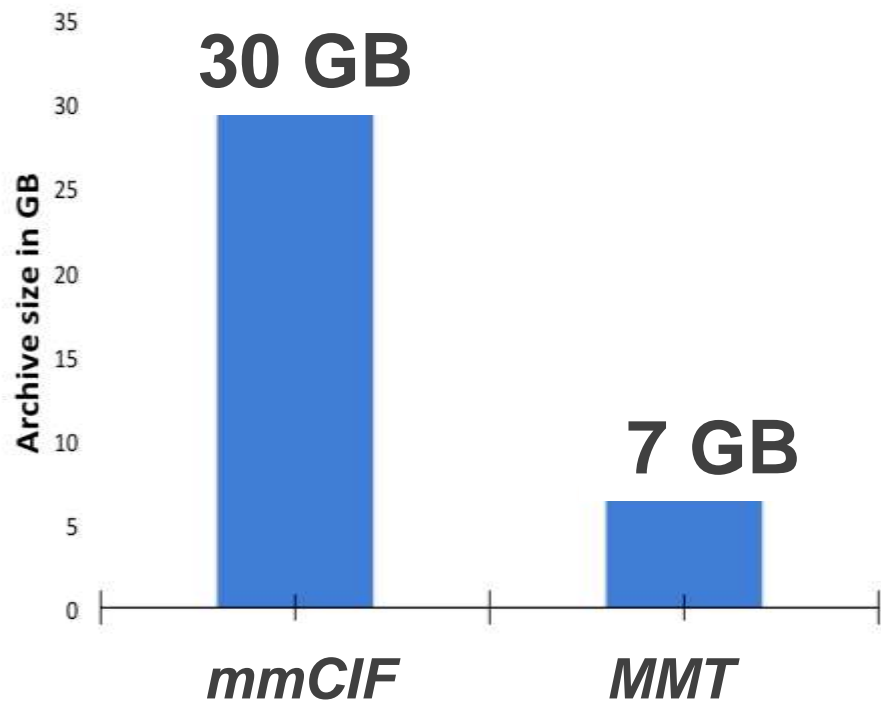
**MessagePack**

*It's like JSON.  
but fast and small.*

# Size and Parsing Speed mmCIF vs. MMTF for 120,000 Structures

Small

Fast

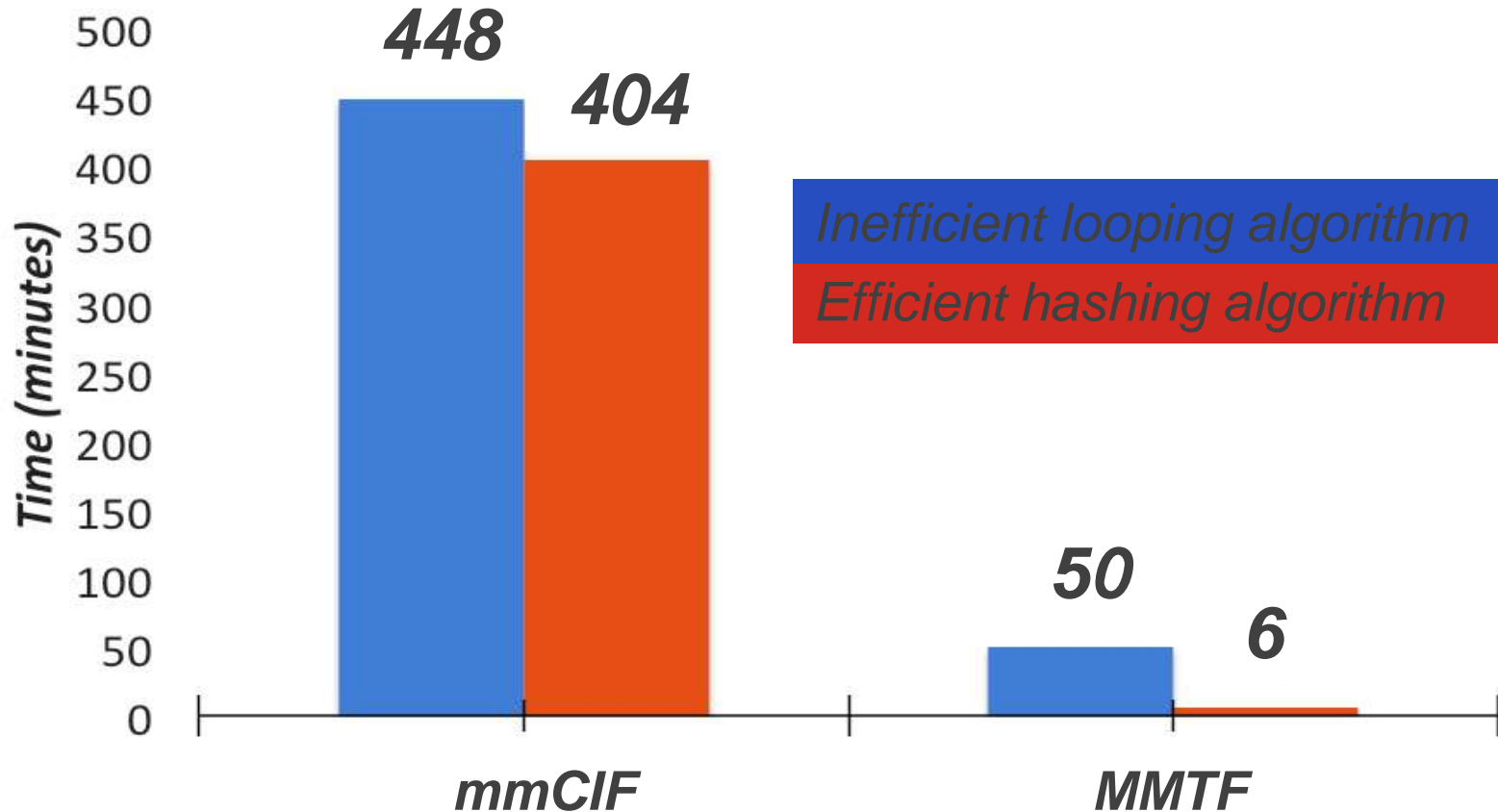


Whole PDB archive GZIP<sup>F</sup> compressed  
(MMTF reduced/lossy: ~800 MB)

Mac mini with 2.6 GHz Intel Core i5<sup>F</sup>  
(4 cores) and 16GB RAM using

# Data Mining using Apache Spark mmCIF vs. MMTF

Find all C-alpha-C-alpha contacts





# Download + Parsing time MMTF vs. mmCIF

Time (seconds) to download\* 100 large PDB structures from UCSD and parse with JavaScript decoder in Chrome browser



\*Note: download times are highly variable and not representative

# Community Engagement

- **Open source specification**
- **Open source decoding libraries**
  - Java
  - JavaScript
  - Python
  - C/C++ (developed by community members)
- **Applications using MMTF**
  - 3Dmol.js, JSmol, iCn3D(NCBI), ICM Viewer, PyMol
  - BioJava, Biopython, MDAnalysis
  - RCSB PDB website

# Summary

- **MacroMolecular Transmission Format (MMTF, [mmtf.rcsb.org](http://mmtf.rcsb.org))**
  - Compressed, binary, efficient representation of 3D structures
    - Lossless representation (~4x compression)
    - Lossy, reduced representation (~37x compression)
- **Compressive Structural Bioinformatics**
  - Algorithms, application, and workflows using MMTF
    - 10 to 100+ fold speedup

## Structure Visualization



## Large Scale PDB Mining



*Web-based molecular graphics for large complexes (2016)*  
*Web 3D '16, 185-186, DOI: 10.1145/2945292.2945324*

# Acknowledgements

Funding: NCI/NIH (U01 CA198942)



*MMTF Early Adopters*



SCHRÖDINGER.

BioJava

NCBI



3Dmol.js

RCSB PDB  
PROTEIN DATA BANK

SDSC SAN DIEGO  
SUPERCOMPUTER CENTER

RCSB PDB

UC San Diego