# Foldcomp: scalable solution for compressing huge protein structure database

**Hyunbin Kim**
**Milot Mirdita**
**Martin Steinegger**

Seoul National University
South Korea

## We compressed AlphaFold Database (23TB) into 1.1TB

### Advent of unprecedentedly large protein structure databases

The AlphaFold databases of 214M UniProt proteins and the ESMatlas catalog of nearly 700M metagenomic protein structures provide valuable resources to the community. However, their extensive sizes of 23TB and 15TB, respectively, exceed the capacity of standard workstations and pose a challenge even to well-equipped cluster environments.

AlphaFold database
214M
23TB tar.gz

AFDB
1.1TB

ESM atlas
700M
15TB pdb.gz in tar.gz

and
more
designed or
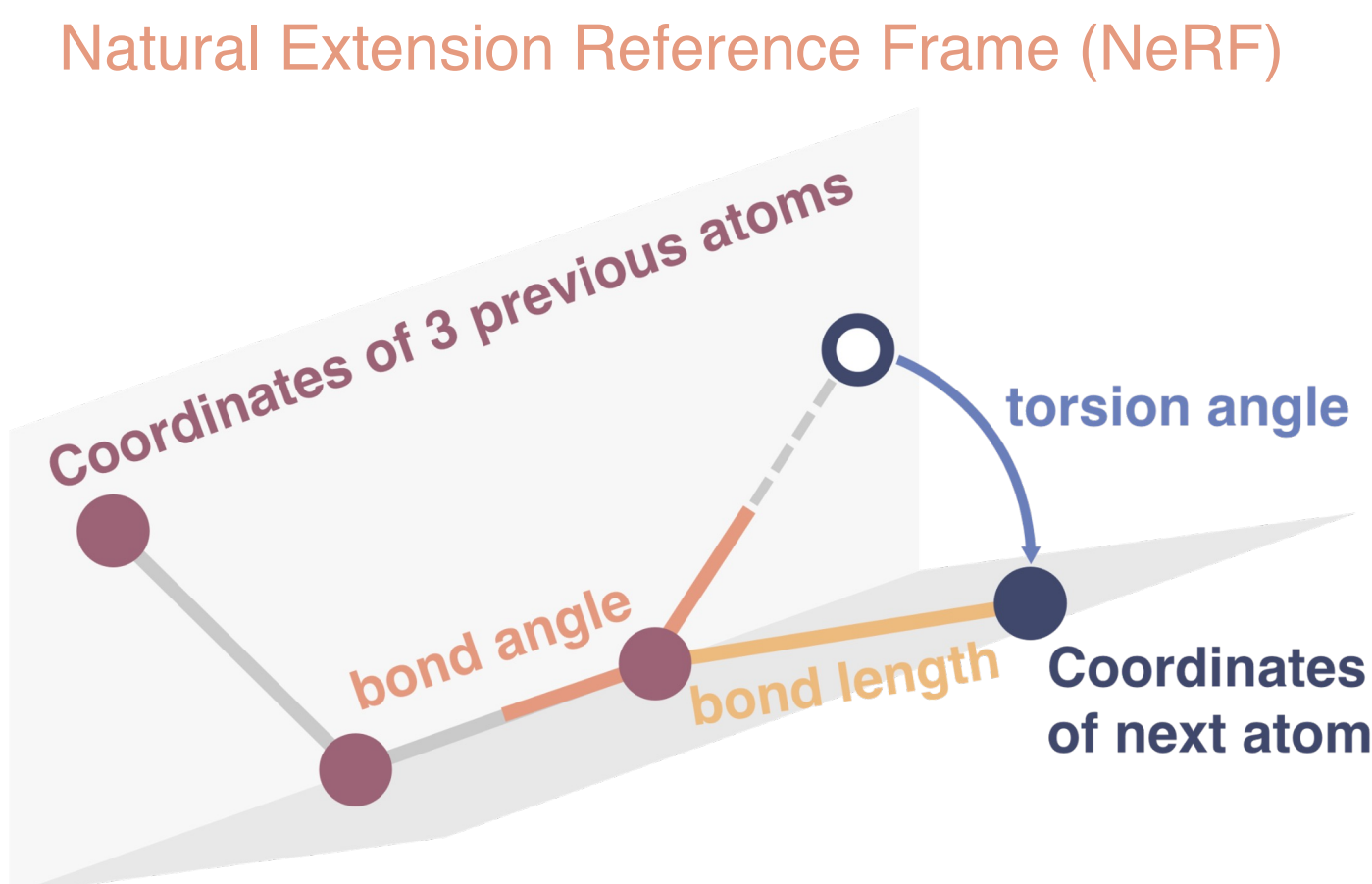metagenomic structures

ESM atlas
1.8TB

PDB

### Novel compression format of torsion angles

We introduce **Foldcomp**, a novel compression algorithm that encodes the torsion and bond angles in a compact binary format, named FCZ. Foldcomp achieves up to 90% compression compared to float-encoded 3D coordinates, requiring only 13 bytes per residue.

**90**% **13**byte/AA **7.7**kb/structure

Compression

**Read PDB/CIF**

**Cartesian to internal**

**Write FCZ**

Architecture

1 byte

Residue 5bit — R

Torsion angles 35bit — omega | psi | phi

Bond angles 24bit — CA-C-N | C-N-CA | N-CA-C

8 bytes per residue

Side chain angles — Asn O | Asn CB | Asn CG | Asn OD1 | Asn ND2 — ~5 bytes per residue

+

Internal anchor coordinates — X | Y | Z — Every 25th AA N, Cα, C — 36 bytes

Internal anchor points to reset error accumulation
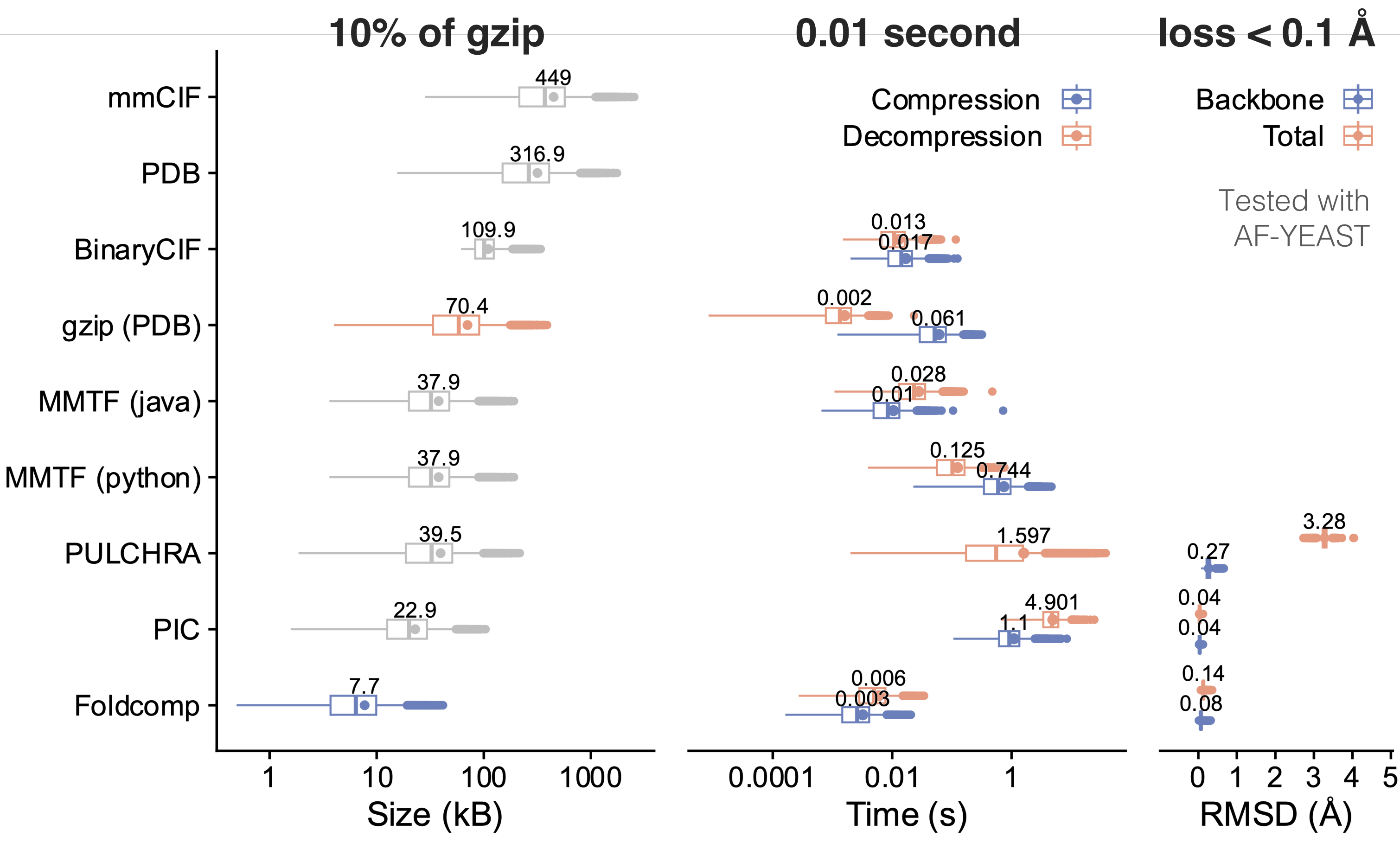
### Reducing loss to less than experimental error

Reconstruction of original coordinates is accomplished by utilizing the NeRF algorithm with internal anchor points. By averaging bi-directional reconstructed coordinates, we were able to reduce reconstruction loss to ~0.08Å range.

Natural Extension Reference Frame (NeRF)

Coordinates of 3 previous atoms
torsion angle
bond angle
bond length
Coordinates of next atom

Decompression

**Read FCZ**

**Restore coordinates**

**Write PDB**

Forward NeRF
Anchor
Reverse NeRF
Weighted average of bi-directional coordinates

### As fast as gzip

Our method is **as fast as gzip**, with 3ms and 6ms for compression and decompression, respectively.

**10**ms **0.08**Å of RMSD

### DB format for less overhead

MMseqs2 database format reduced unnecessary padding bytes and overhead from file numbers

Loose files (.fcz)
FCZs

Foldcomp database

| 1 FCZ | | | 1 | 0 | 1612 | | 1 | AF-Q0... |
| ⋮ | | | ⋮ | | | | ⋮ | |
| N FCZ | | | N | 44849 | | | N | AF-PZ... |

FCZ entries — Index — Lookup — start & size — name

### Foldcomp is publicly available

as a command line interface and a Python API at https://foldcomp.foldseek.com.

## pip install foldcomp

Foldcomp is compatible with **Foldseek** and additionally, has been augmented by community contributions, such as a **PyMol** plugin and a dataset wrapper in **Graphein**. We provide the compressed database of AlphaFold database (1.1TB), ESMatlas (1.8TB), SwissProt (2.9GB), and recently released AlphaFold2 cluster representatives (2.2GB) at https://foldcomp.steineggerlab.workers.dev.

**10% of gzip** — **0.01 second** — **loss < 0.1 Å**

Compression / Decompression
Backbone / Total
Tested with AF-YEAST

| | Size (kB) | Time (s) | RMSD (Å) |
|---|---|---|---|
| mmCIF | 449 | | |
| PDB | 316.9 | | |
| BinaryCIF | 109.9 | 0.013 / 0.017 | |
| gzip (PDB) | 70.4 | 0.002 / 0.061 | |
| MMTF (java) | 37.9 | 0.028 / 0.01 | |
| MMTF (python) | 37.9 | 0.125 / 0.744 | |
| PULCHRA | 39.5 | 1.597 | 3.28 |
| PIC | 22.9 | 4.901 / 1.1 | 0.27 / 0.04 / 0.04 |
| Foldcomp | 7.7 | 0.006 / 0.003 | 0.14 / 0.08 |

## References

AlphaFold database – Varadi et al., *Nucleic acids research* (2022).
ESMatlas – Lin et al., *Science* (2023).
Foldseek – van Kempen et al., *Nature biotechnology* (2023).
MMseqs2 – Steinegger and Söding., *Nature biotechnology* (2017).
NeRF – Parsons et al., *Journal of computational chemistry* (2005).
Foldcomp is published at https://doi.org/10.1093/bioinformatics/btad153.

NRF

Foldcomp Github repository

Foldcomp Databases