# Fast lossy protein structure compression algorithm

Hyunbin Kim[1], Johannes Söding[2], Martin Steinegger[1]

[1]Seoul National University, South Korea; [2]Max Planck Institute, Germany; martin.steinegger@snu.ac.kr

## ABSTRACT

AlphaFold2 produces structure predictions at high quality and speed. EMBL and DeepMind have announced to soon release a database containing over 100 million predicted structures covering the UniRef90. Thus, a future with billions of predicted structures is soon imaginable. Additionally, the prediction speed is constantly improved. e.g., ColabFold is ~100x faster compared to baseline AF2.

However, with advances in speed, storing all the structures is becoming a major issue. Storing the structure of a protein with 250 residues in PDB format takes approx. 200 kilobytes (only 3D coordinates 25 kilobytes), thus one billion structures would require hundreds of terabytes.
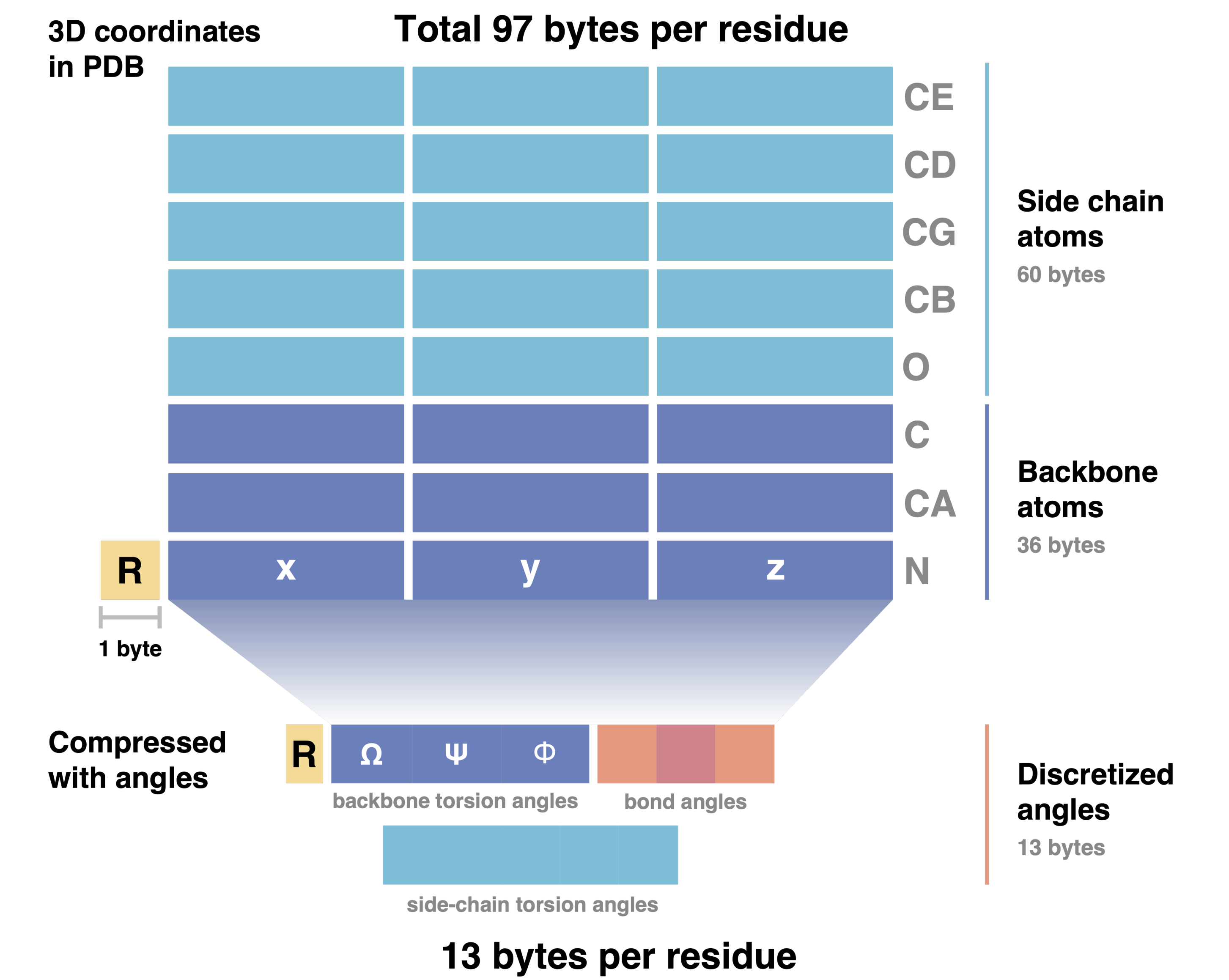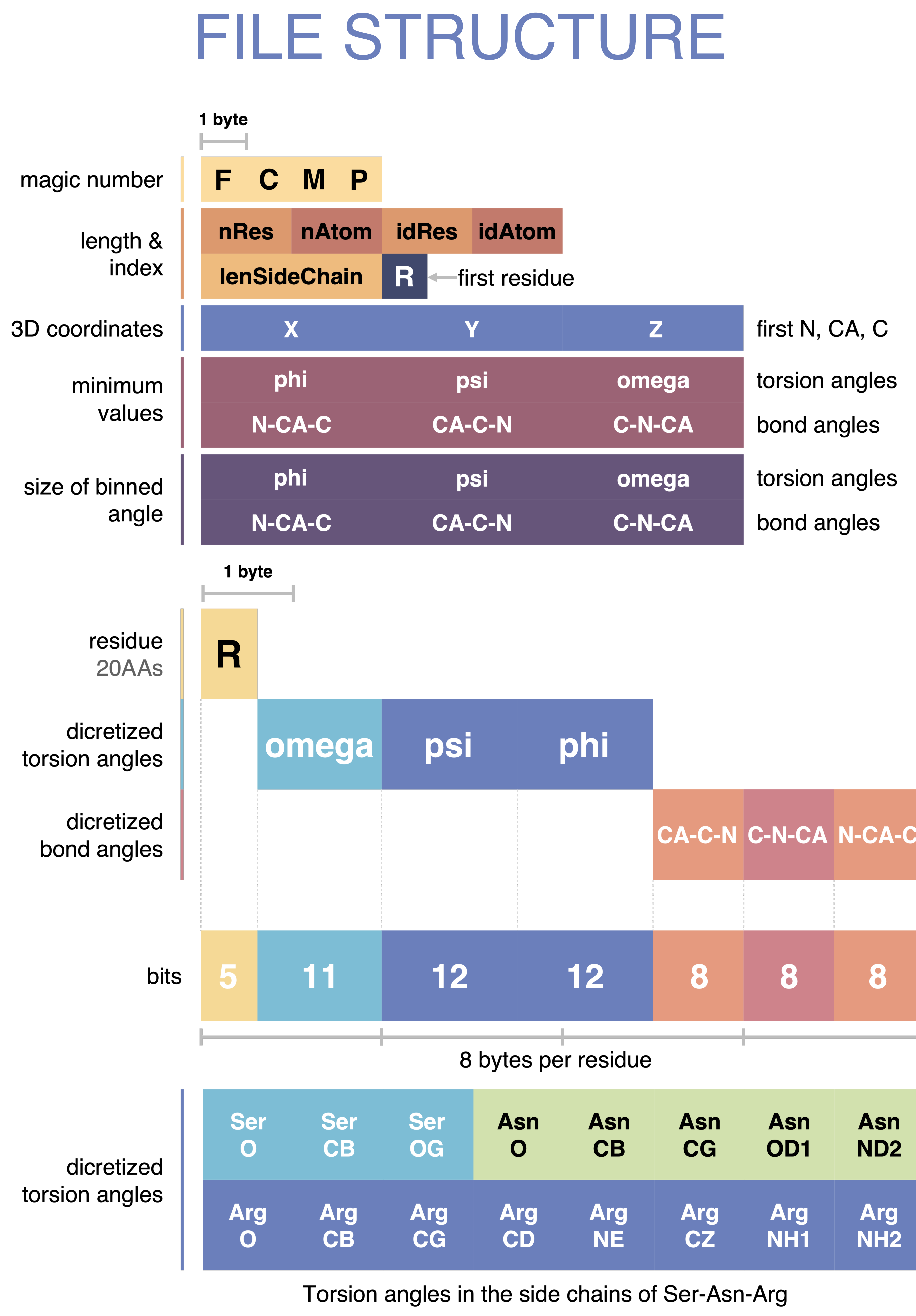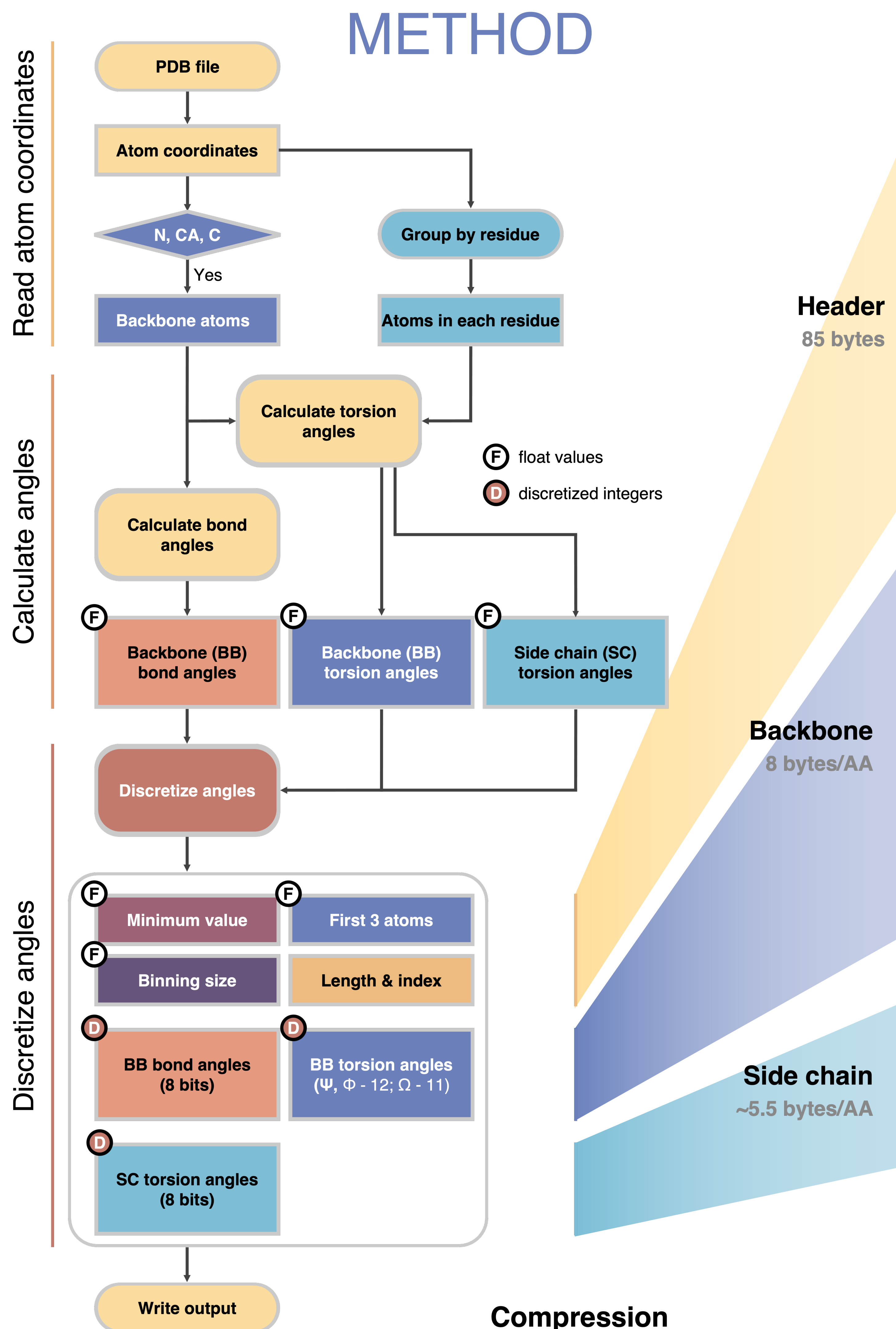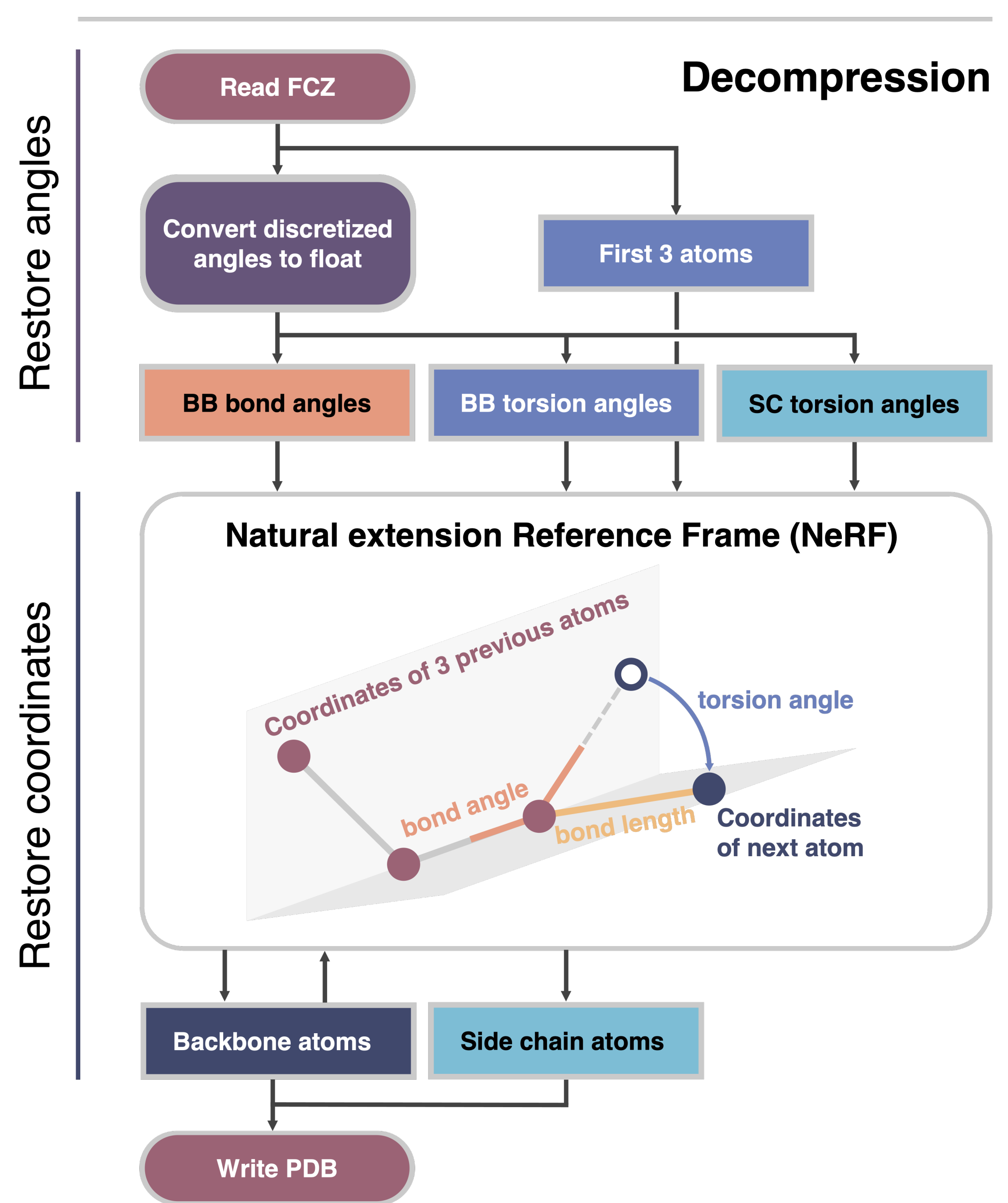
Here, we propose a novel format and method to compress protein structures requiring only 10 kilobytes for a protein structure of average size (4.8 kb for coordinates), reducing the required storage space by an order of magnitude.

## METHOD



Read atom coordinates / Calculate angles / Discretize angles

PDB file → Atom coordinates → N, CA, C → Backbone atoms → Calculate torsion angles → Calculate bond angles

F float values
D discretized integers

Backbone (BB) bond angles / Backbone (BB) torsion angles / Side chain (SC) torsion angles → Discretize angles

Minimum value / First 3 atoms / Binning size / Length & index
BB bond angles (8 bits) / BB torsion angles ($\Psi$, $\Phi$ - 12; $\Omega$ - 11)
SC torsion angles (8 bits) → Write output

**Compression** / **Decompression**

Read FCZ → Convert discretized angles to float / First 3 atoms → BB bond angles / BB torsion angles / SC torsion angles

**Natural extension Reference Frame (NeRF)**
Coordinates of 3 previous atoms, torsion angle, bond angle, bond length, Coordinates of next atom

Backbone atoms / Side chain atoms → Write PDB

## FILE STRUCTURE



**Header** 85 bytes — magic number: F C M P; length & index: nRes, nAtom, idRes, idAtom, lenSideChain, R first residue; 3D coordinates: X Y Z first N, CA, C; minimum values: phi, psi, omega (torsion angles), N-CA-C, CA-C-N, C-N-CA (bond angles); size of binned angle: phi, psi, omega (torsion angles), N-CA-C, CA-C-N, C-N-CA (bond angles)

**Backbone** 8 bytes/AA — residue 20AAs: R; dicretized torsion angles: omega, psi, phi; dicretized bond angles: CA-C-N, C-N-CA, N-CA-C; bits: 5, 11, 12, 12, 8, 8, 8 — 8 bytes per residue

**Side chain** ~5.5 bytes/AA — dicretized torsion angles: Ser O, Ser CB, Ser OG, Asn O, Asn CB, Asn CG, Asn OD1, Asn ND2 / Arg O, Arg CB, Arg CG, Arg CD, Arg NE, Arg CZ, Arg NH1, Arg NH2

Torsion angles in the side chains of Ser-Asn-Arg

**3D coordinates in PDB** — Total 97 bytes per residue

Side chain atoms 60 bytes: CE, CD, CG, CB, O
Backbone atoms 36 bytes: C, CA, N
R x y z — 1 byte

**Compressed with angles** — R $\Omega$ $\Psi$ $\Phi$ (backbone torsion angles), bond angles, side-chain torsion angles — Discretized angles 13 bytes

**13 bytes per residue**

We achieve this reduction by efficiently encoding the torsion angles of the backbone as well as the side-chain angles in a compact format. We show that using our lossy compression has no impact on structural downstream analysis. By storing angles with an optimized bit-format, we can reduce the storage required by 90% compared to float-encoded 3D coordinates, while maintaining a high compression and decompression speed.

## BENCHMARK RESULT



AF-Yeast 6,040 structures
gzip -6 / pulchra / PIC / foldcomp (ours)
Compressed file size avg. **10.2Kb**
Running time avg. **0.4s**

Loss in compression

| PDB ID | Tool | RMSD |
|---|---|---|
| 1a0fA | foldcomp | **0.227** |
| | pic | 18.976 |
| | pulchra | 3.208 |
| 1a0aA | foldcomp | **0.154** |
| | pic | 20.688 |
| | pulchra | 3.343 |
| 1a0p_ | foldcomp | 6.744 |
| | pic | 19.091 |
| | pulchra | **3.476** |
| 1a0i_ | foldcomp | 4.241 |
| | pic | 21.420 |
| | pulchra | **3.370** |
| 1a0tP | foldcomp | **0.443** |
| | pic | 20.958 |
| | pulchra | 3.120 |

5 randomly selected PDB files