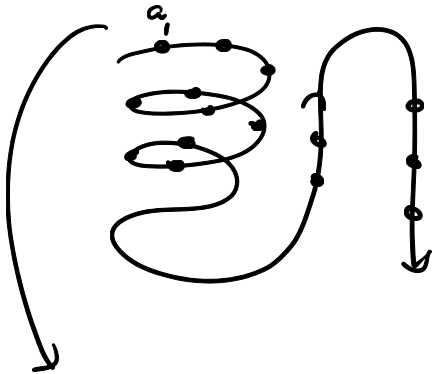
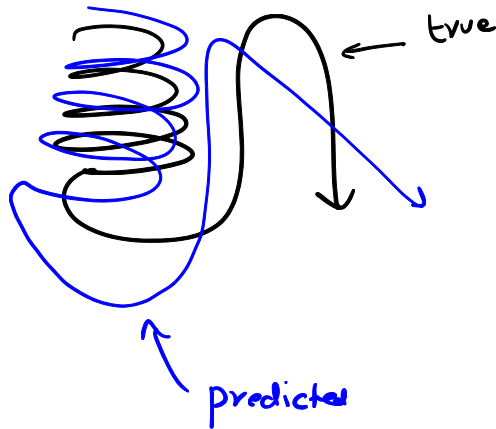


Protein Structure Prediction



$\alpha_1 \rightarrow (x_1, y_1, z_1)$ 3D coords vs $(\phi_1, \psi_1, \omega_1)$ torsion angles
 amino acid $G \leftarrow$ primary structure (seq)



$$\sum_{i=1}^N (\vec{T}_i - R \vec{P}_i)^2 = \text{Obj: } \min_R$$

N : length of protein

\vec{T}_i : (x_i, y_i, z_i) true

\vec{P}_i : (x_i, y_i, z_i) predicted

R : 3×3 rotation matrix

\nwarrow the best possible rotation matrix

$$\vec{\mu}_T = \frac{1}{N} \sum \vec{T}_i = \frac{1}{N} \sum \vec{P}_i = \vec{\mu}_P$$

means are super-imposed first

CRMSD : RMSE root mean squared error

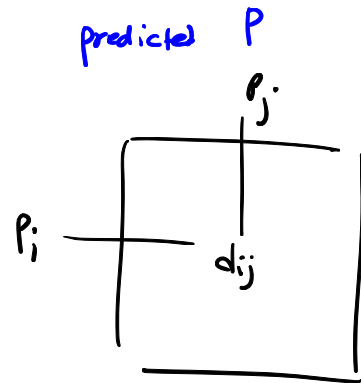
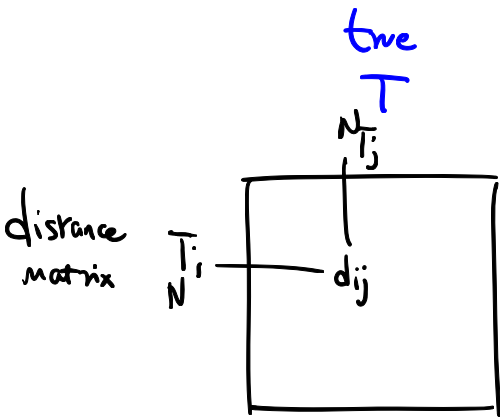
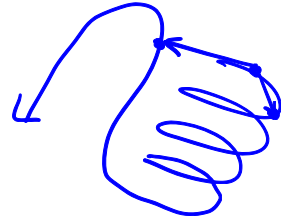
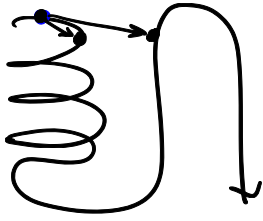
$$\vec{P}'_i = R \vec{P}_i$$

Cost

$$\sqrt{\frac{1}{N} \sum (\vec{T}_i - \vec{P}'_i)^2}$$

$$\vec{T}_i, \vec{P}_i, \vec{P}'_i \in \mathbb{R}^3$$

distance RMSD



$$d_{ij} = \|\vec{T}_i - \vec{T}_j\|$$
$$= \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2}$$

D_p : distance matrix for P

D_T : Symmetric $N \times N$ distance matrix for T

$$\|D_T - D_P\|_F = \text{dRMSD} = \sqrt{\frac{1}{N^2} \sum_i \sum_j (D_T(i,j) - D_P(i,j))^2}$$

Frobenius Norm

Simplified Distance Matrix : Contact Map

$$\theta = 8\text{\AA}^\circ$$

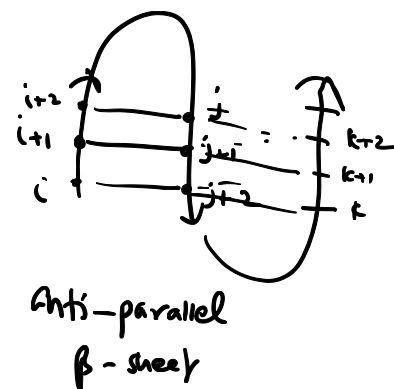
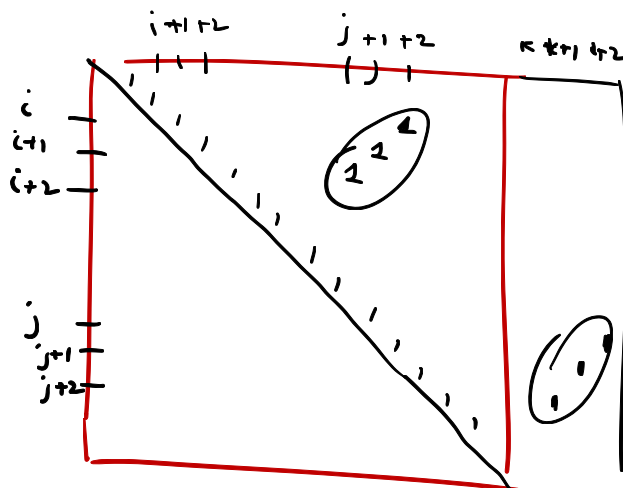
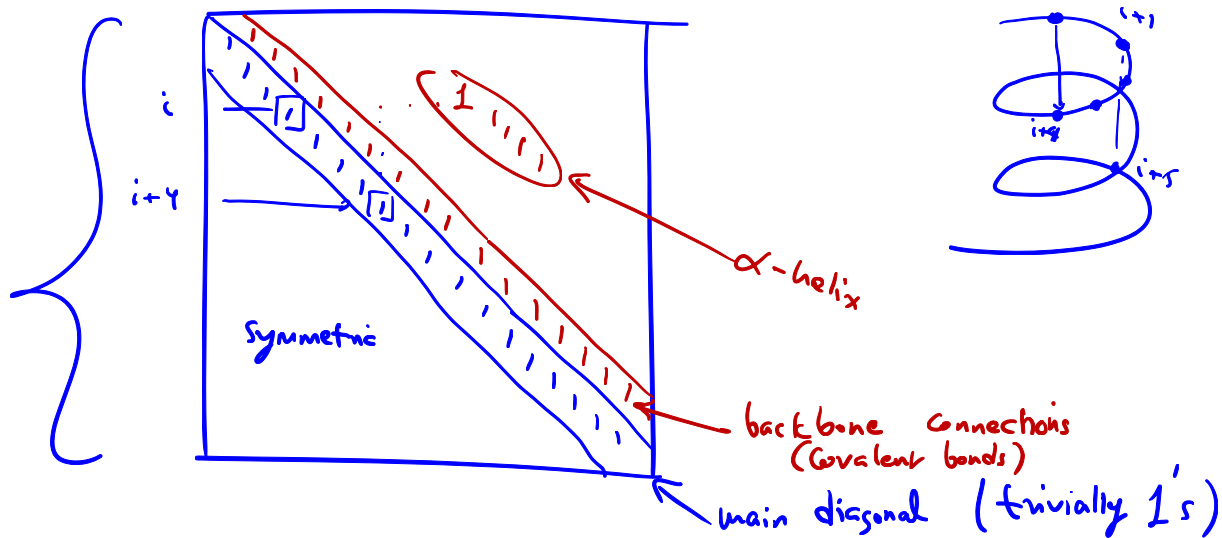
$$A^\circ = 10^{-10} \text{ m} = 0.1 \text{ nm}$$

	τ_1	τ_2	τ_3
τ_1	0	6\AA°	10\AA°
τ_2	6\AA°	0	6\AA°
τ_3	10\AA°	6\AA°	0

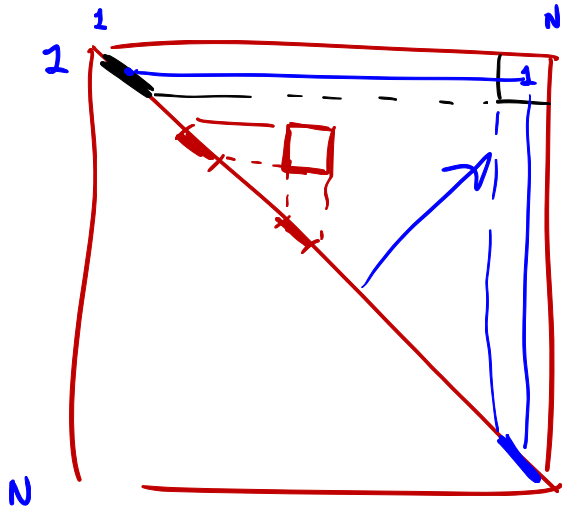
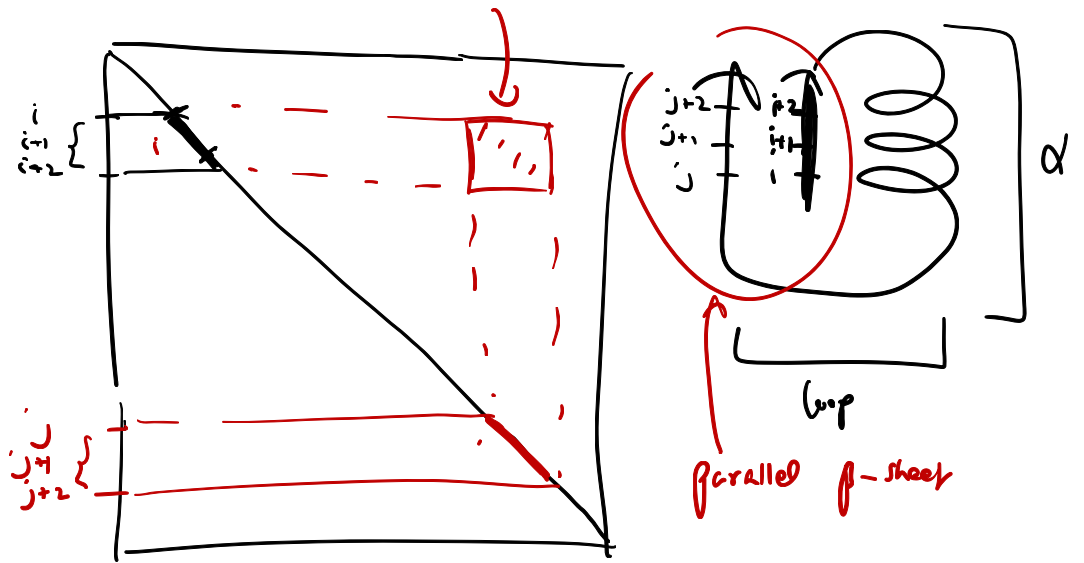


1	1	0
1	1	1
0	1	1

N

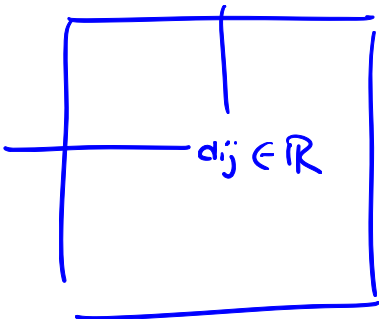


off-diagonal

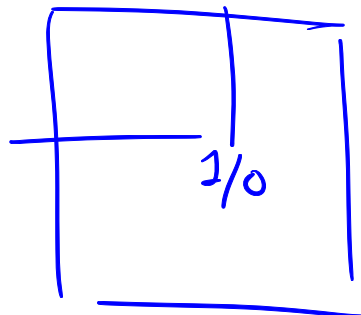


Very sparse
 $O(N)$ contacts

Distance map

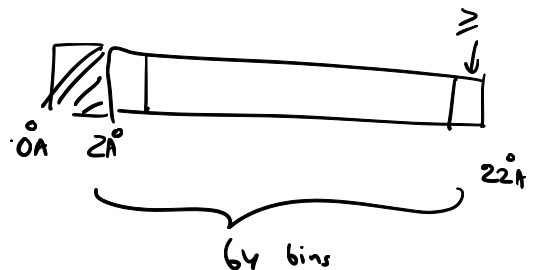
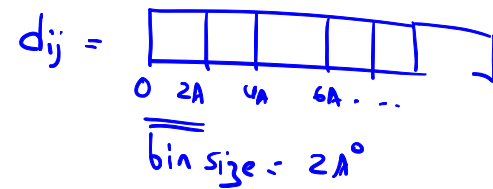


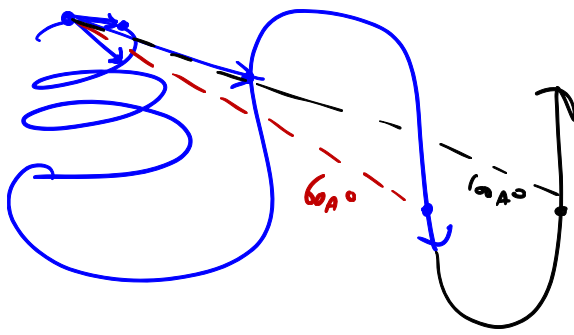
Contact map (Θ)



bins = 2
 $\Theta = 8A^\circ$

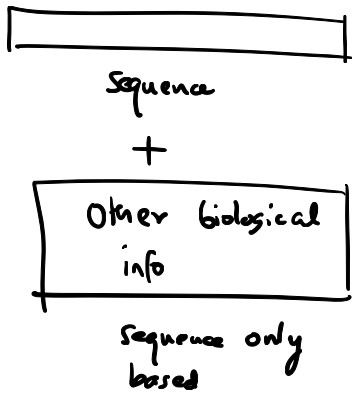
discretized distance map





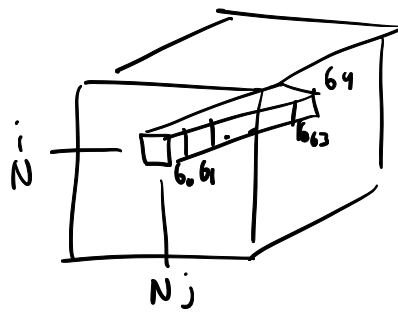
b_0, b_1, \dots, b_{63}
 \uparrow
 $\geq 22\text{\AA}$

Input

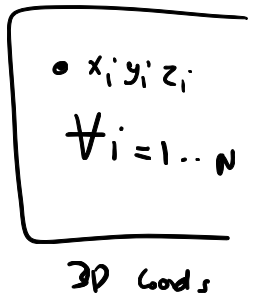


Output

64-bin distance map



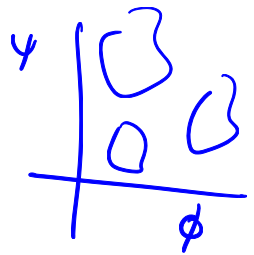
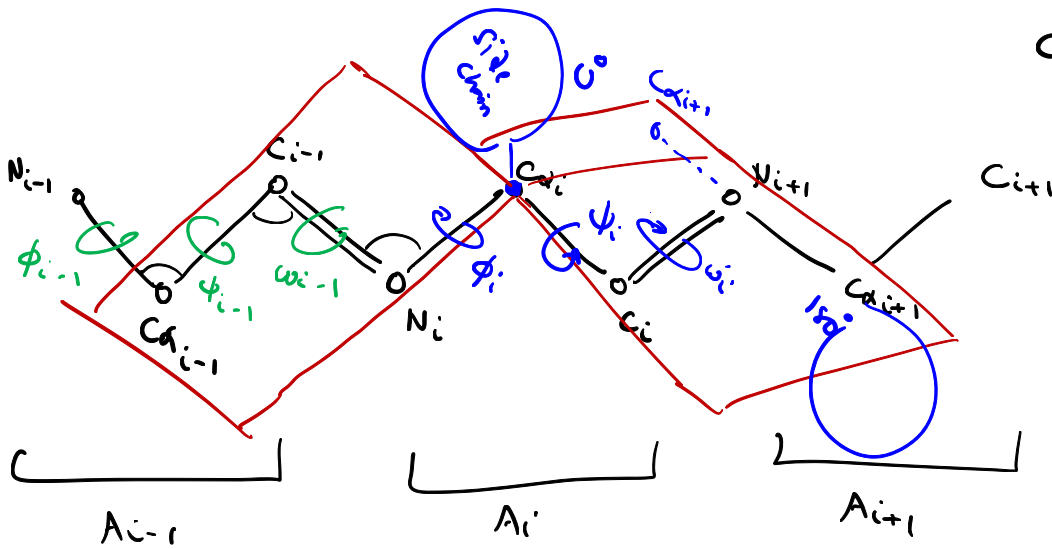
Final output



Torsion Angles (dihedral angle)

N-terminus

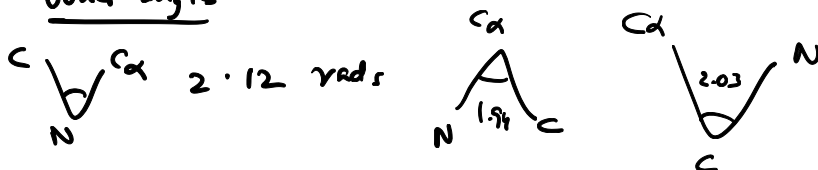
C-terminus



bond-length

- C-N : 1.33\AA
- N-C α : 1.46\AA
- C α -C : 1.52\AA

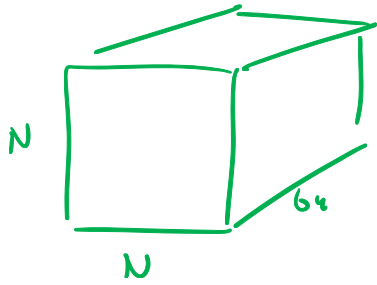
bond angle



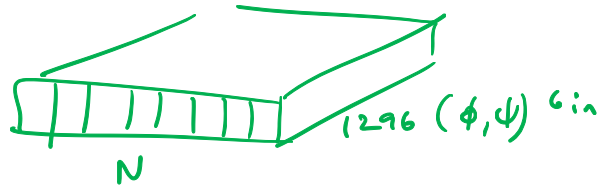
$\approx 110^\circ - 115^\circ$

$w_i \in \{0^\circ, 180^\circ\}$
 \uparrow \uparrow
 cis trans
 \checkmark
 pre-dominant one

$w_i = 180^\circ$ in practice
 with rare
 exceptions!



+



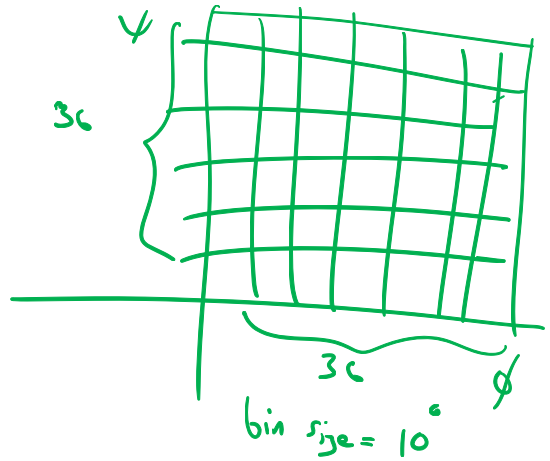
discrete distance map

torsion angle bins

$(CE_d + CE_a)$

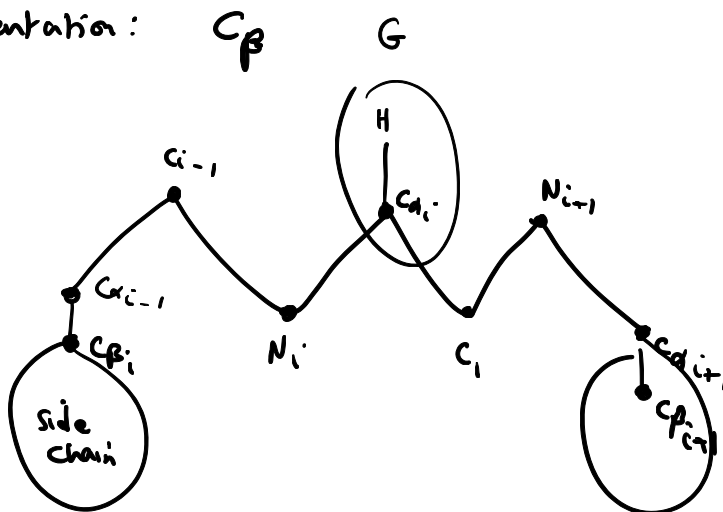
final loss

(ASA: accessible surface area, etc.)



1296 bins

Representation:



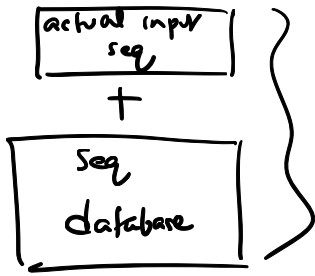
$C_{\beta i} \forall i = 1 \dots N$
 except for G
 which uses C_α

$$d_{ij} = \left\| \begin{matrix} \vec{C}_{\beta i-1} \\ (x \ y \ z) \end{matrix} - \begin{matrix} \vec{C}_{\beta i} \\ (x \ y \ z) \end{matrix} \right\|$$

AlphaFold 1 : Very deep CNN

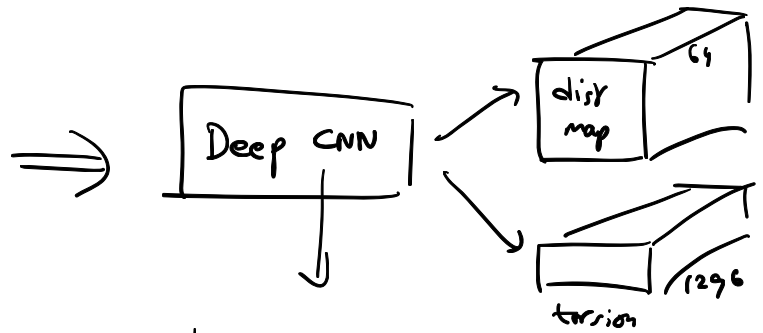
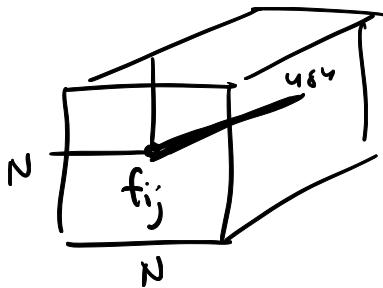
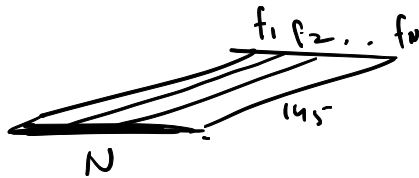
$|N| = \text{length}$

MILGRAMS...



f_i : feature vector for pos i in seq (N)
 $\in \mathbb{R}^{145}$

f_{ij} : pair-wise features (N^2)
 for pos i & pos j
 $\in \mathbb{R}^{484}$



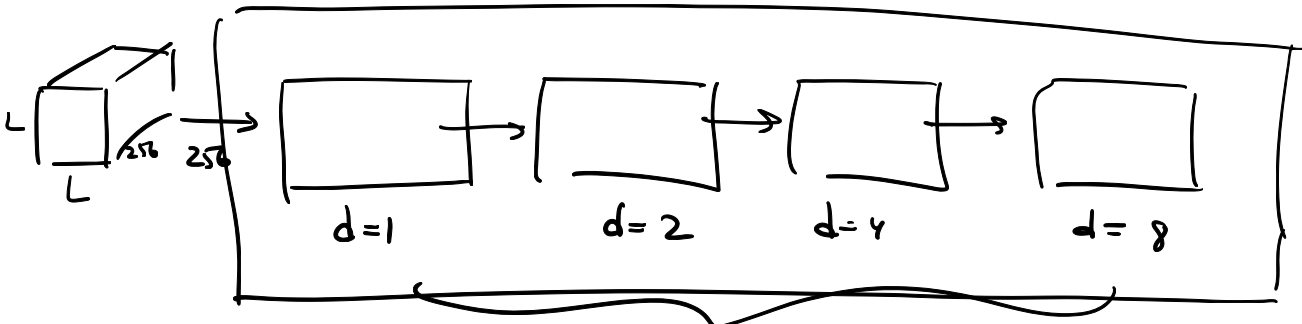
dilations

dilated convolutions

220 blocks of dilated CNNs

7x7 blocks
 of dim 256

48x48 blocks
 of 128 dim features



$d \equiv \text{dilation}$

$L: 64 \Rightarrow 64 \times 64$ crop of the distance matrix!

220 x 6 layers = 1320 layers!