

AlphaFold 1

Q: Protein sequence, length = L
MLIGAWS...

Homologous sequence
related by evolution

PSSM: position specific
Scoring matrix

MSA: Multiple sequence alignment

- 1) query the seq DB
- 2) extract similar sequence
- 3) MSA
- 4) PSSM

BLAST

Q → MLIGAWSPI M
→ MLGAPISTW
→ M GAWP
→ MLIGSWSTAW

CASP 15
Critical Assessment of
Structure Prediction
(2022!)

FM
Free Modeling

TBM
Template-based
Modeling

PDB
Protein data Bank
~ 200k structures

Uniprot sequence database
10⁶ + sequences

MSA
↓
weighted version of
LCS
(longest common
subsequence)

Q

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| M | L | I | G | A | W | - | S | P | I | M | - |
| M | L | - | G | A | P | I | S | T | W | - | - |
| M | - | - | G | A | W | - | - | P | - | - | - |
| M | L | I | G | S | W | S | - | T | A | W | - |

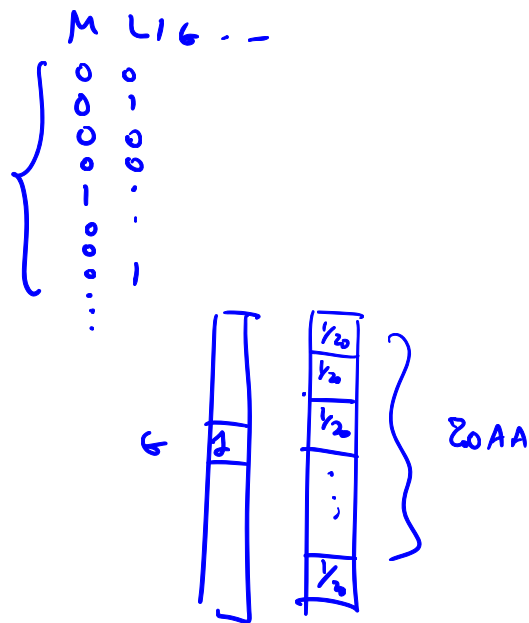
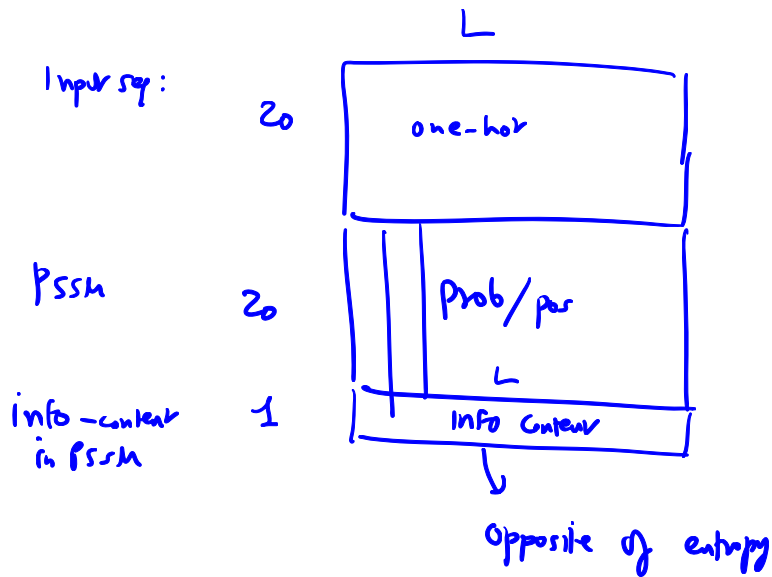
MSA ≡



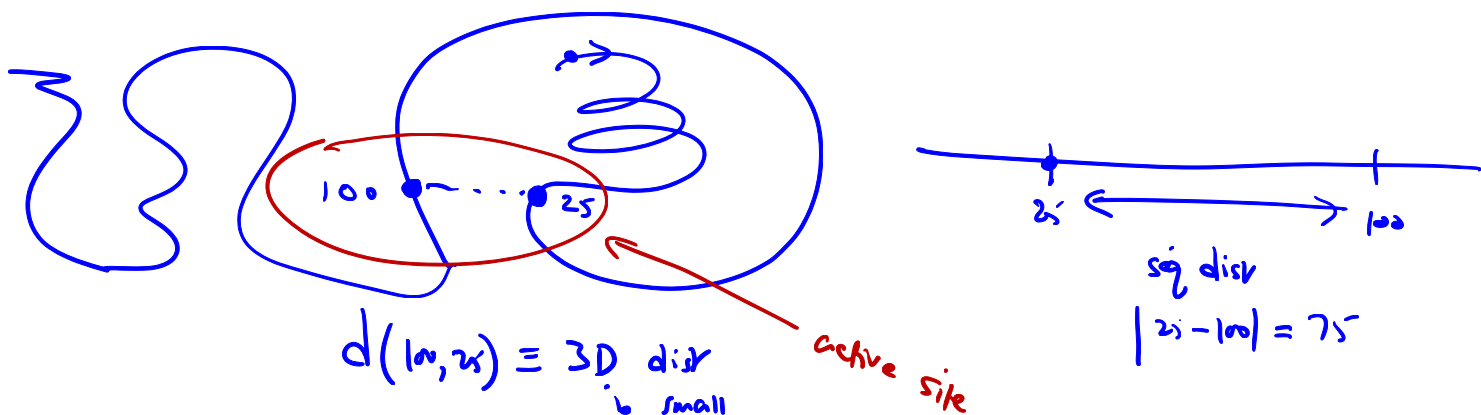
Profile
PSSM

| | | | | | | | | | | |
|----|------|---|------|------|---|---|---|---|---|---|
| | A | L | I | G | W | S | P | I | M | - |
| 20 | 0.75 | | 0.5 | 0.25 | | | | | | |
| 21 | | | | 1.0 | | | | | | |
| L | | | 0.25 | | | | | | | |

Ungapped PSSM : disallow gaps



Derive pair-wise feature from seq



Co-evolution info (Potts model)

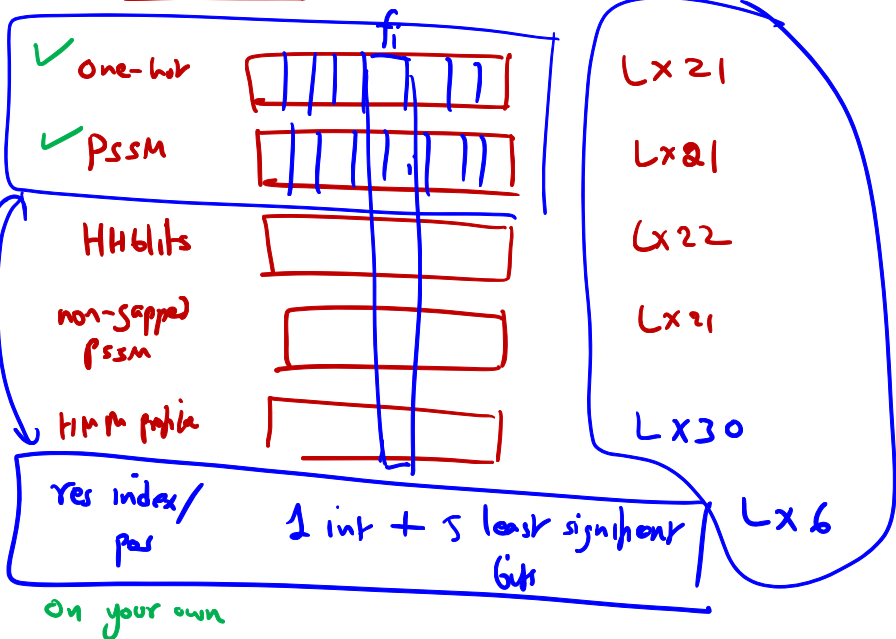
$a_i \text{ vs } a_j \rightarrow f_{ij}$

22 \times 22 = 484 + 1 = 485 dim feature vector

info-content

$f_{ij} \in \mathbb{R}^{485}$

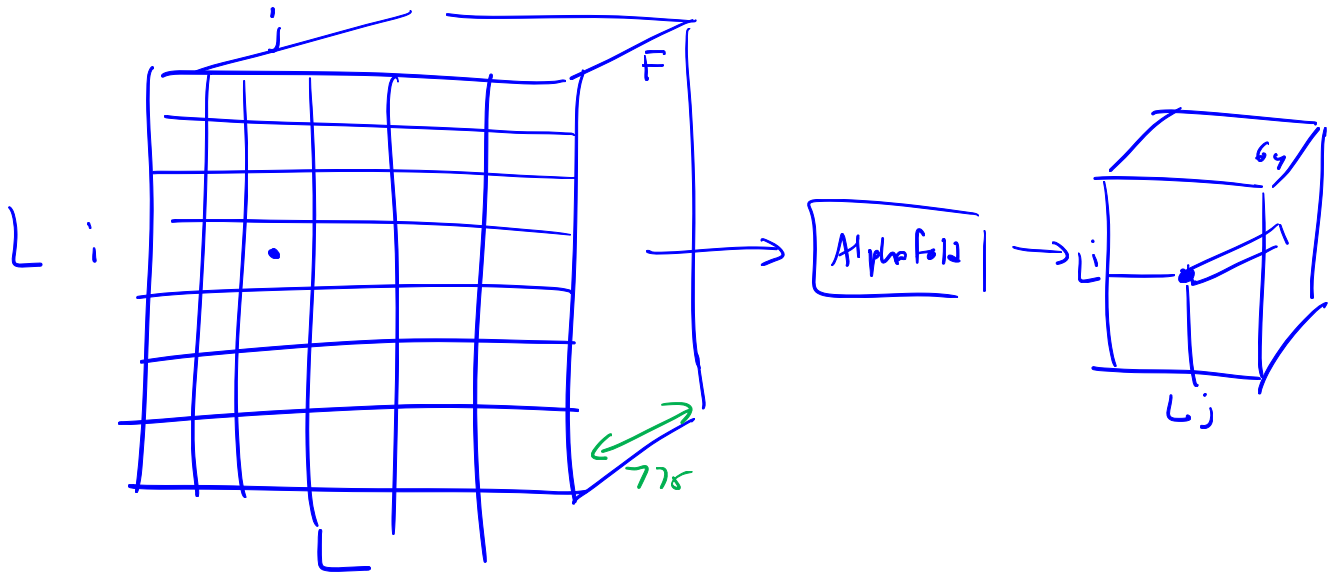
Alpha-fold



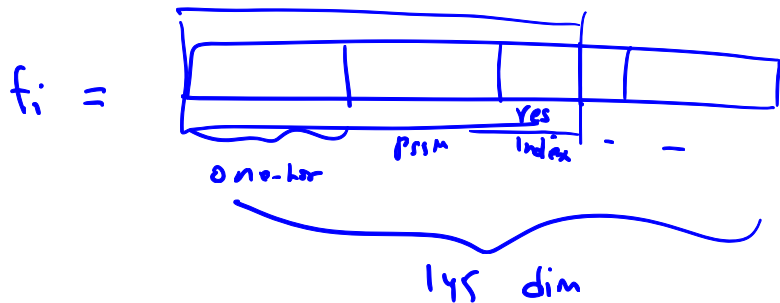
$$f_i \in \mathbb{R}^{145} \quad \forall i = 1 \dots L$$

$$f_{ij} \in \mathbb{R}^{495} \quad \forall i, j = 1 \dots L$$

True distance map ✓



$F_{ij} \equiv$ feature vector for (i, j) cell. extra!
 \equiv \vec{f}_{ij} $\parallel \vec{f}_i \parallel \vec{f}_j$ $\parallel \|\vec{f}_i - \vec{f}_j\| \parallel \vec{f}_i \otimes \vec{f}_j \parallel$
 Co-evolution vector $145 \quad 145 \quad 145 \quad 145$
 495

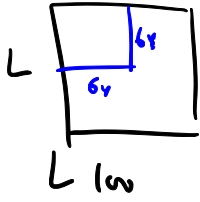


$$\begin{array}{r}
 145 \\
 145 \\
 145 \\
 \hline
 435
 \end{array}$$

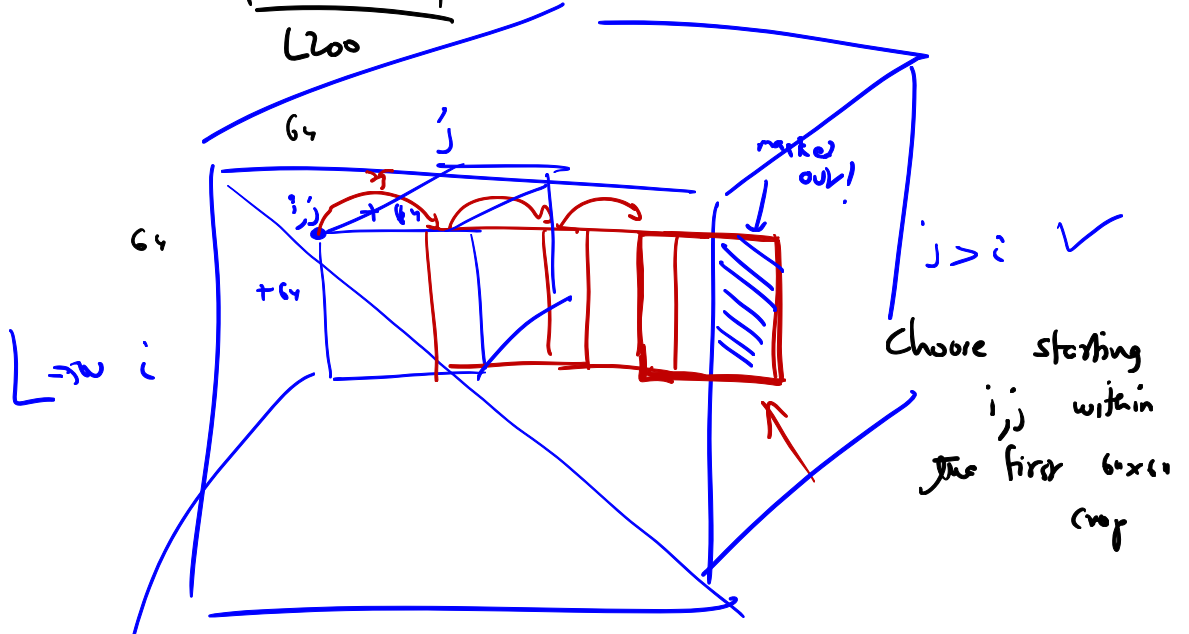
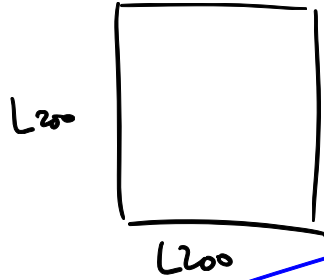
dim features!

Alpha fold

P_1

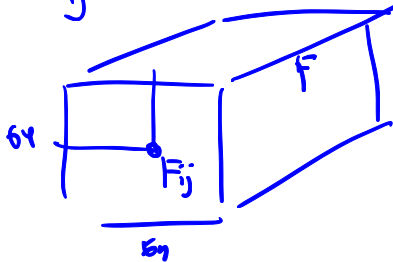


P_2

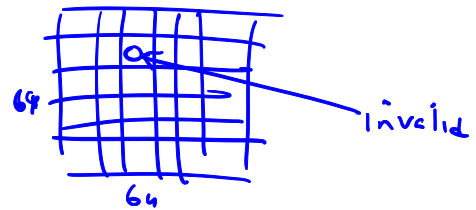


$L = 500$

Use padding of 32 all around

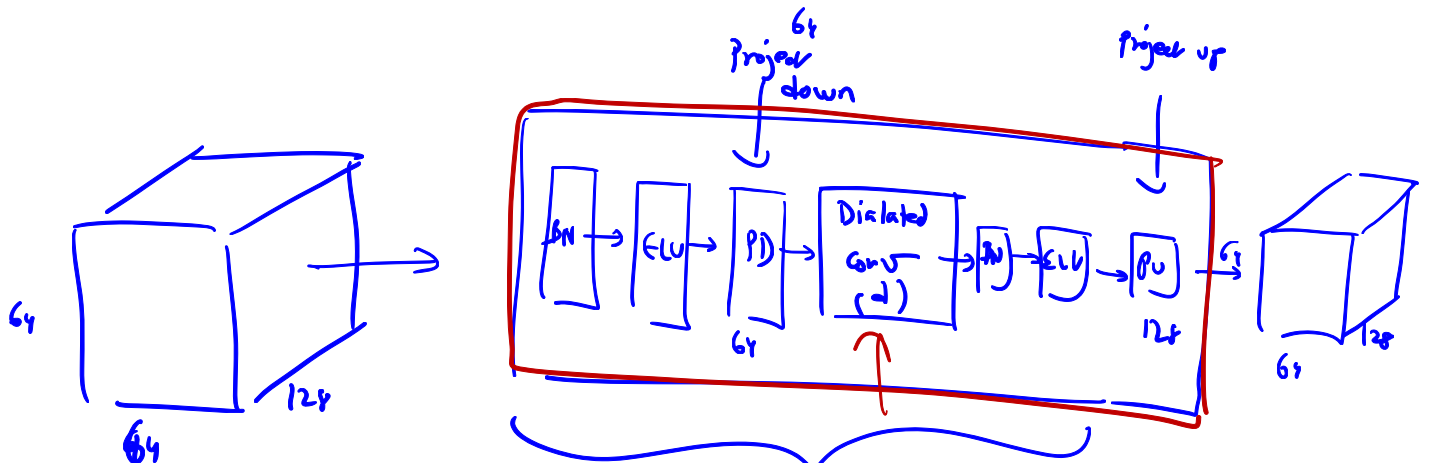
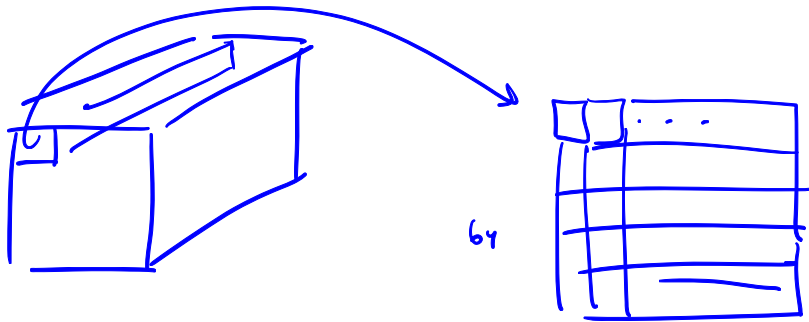
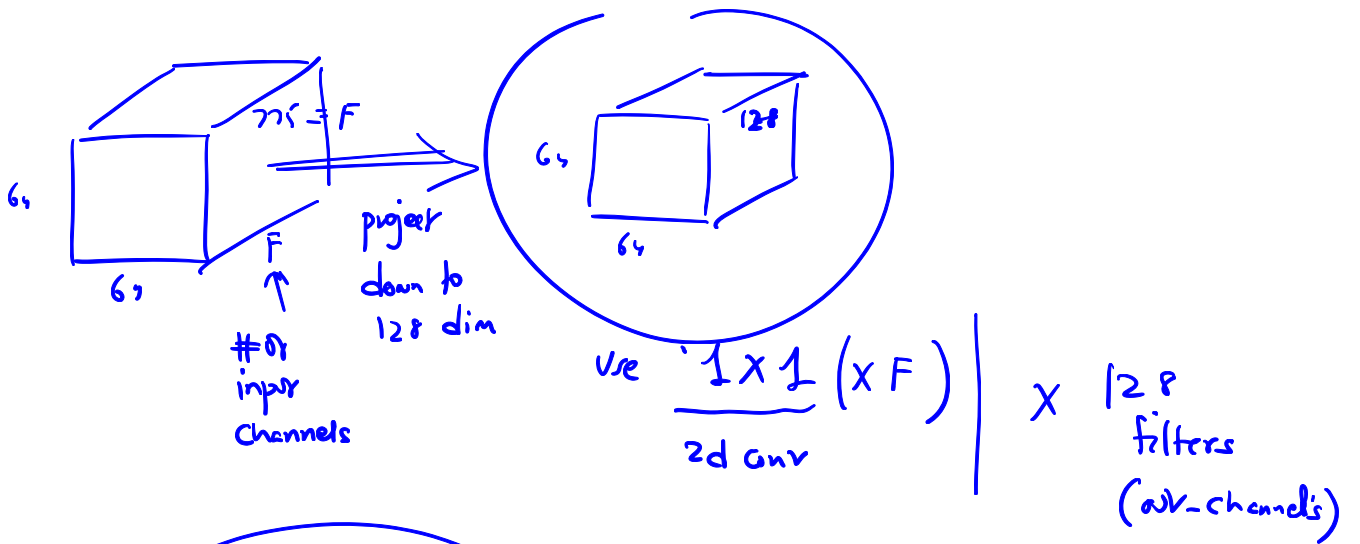


mask

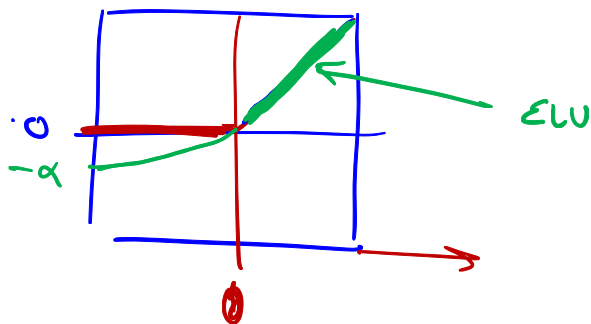


for some protein

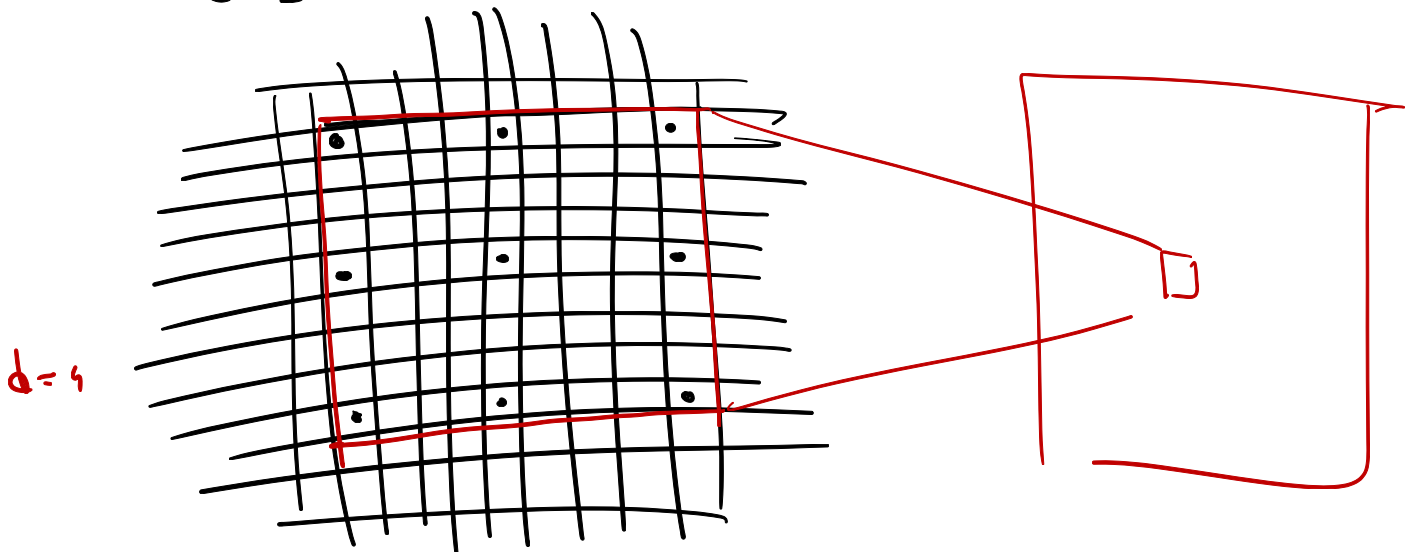
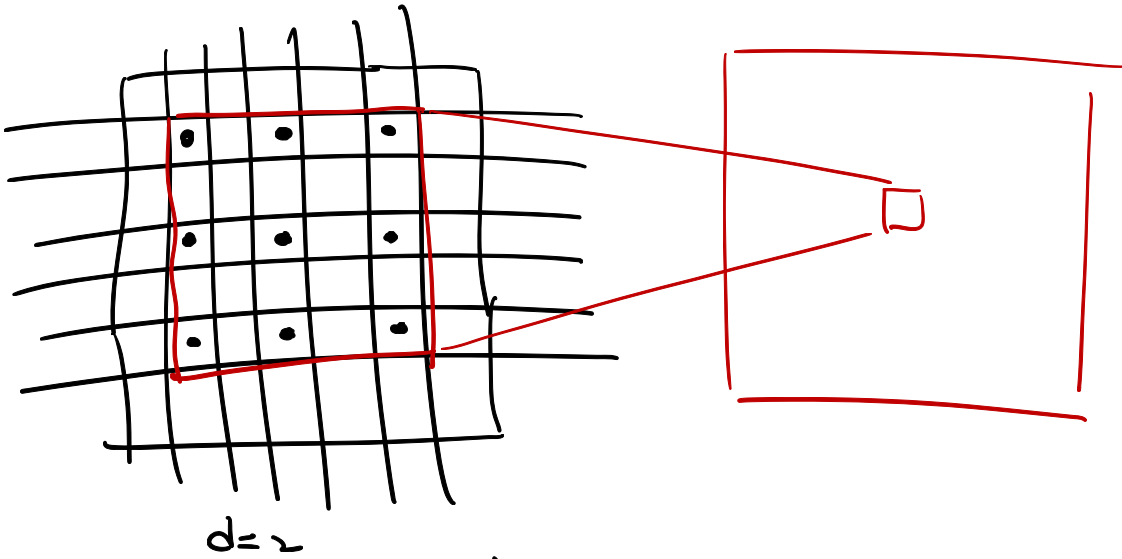
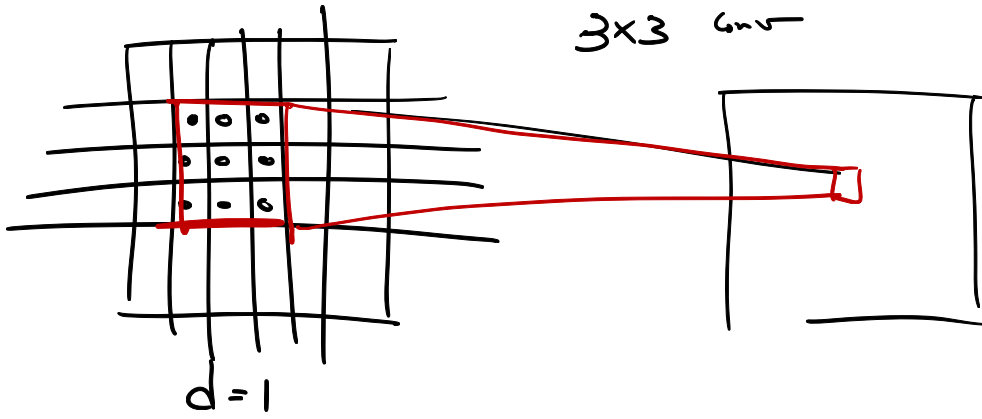
↳ multiple crops with i, j starting + slide

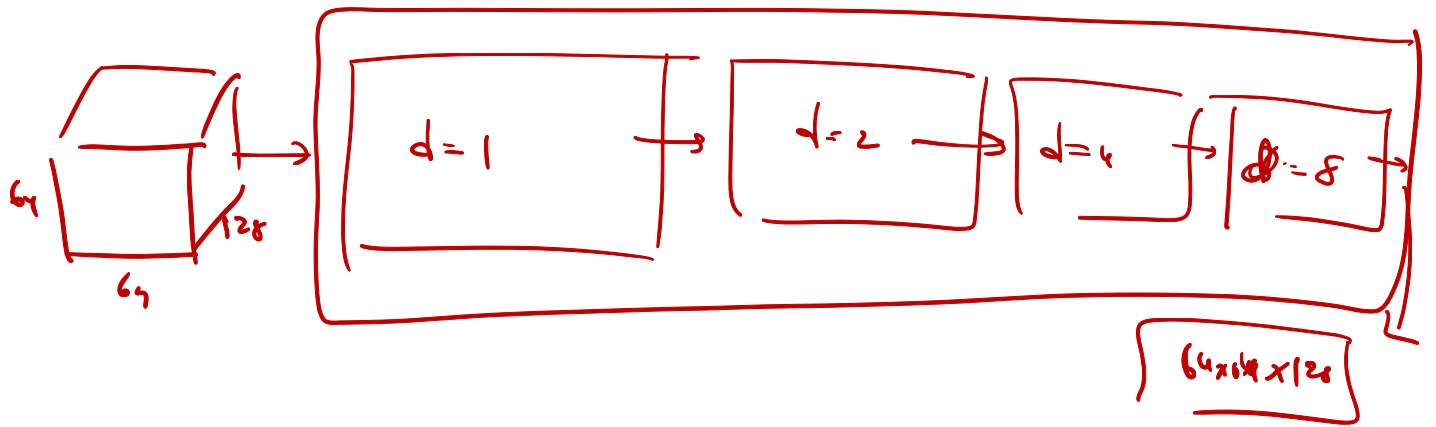


PD/PU
: 1×1 conv
(out-channels
in-channels)

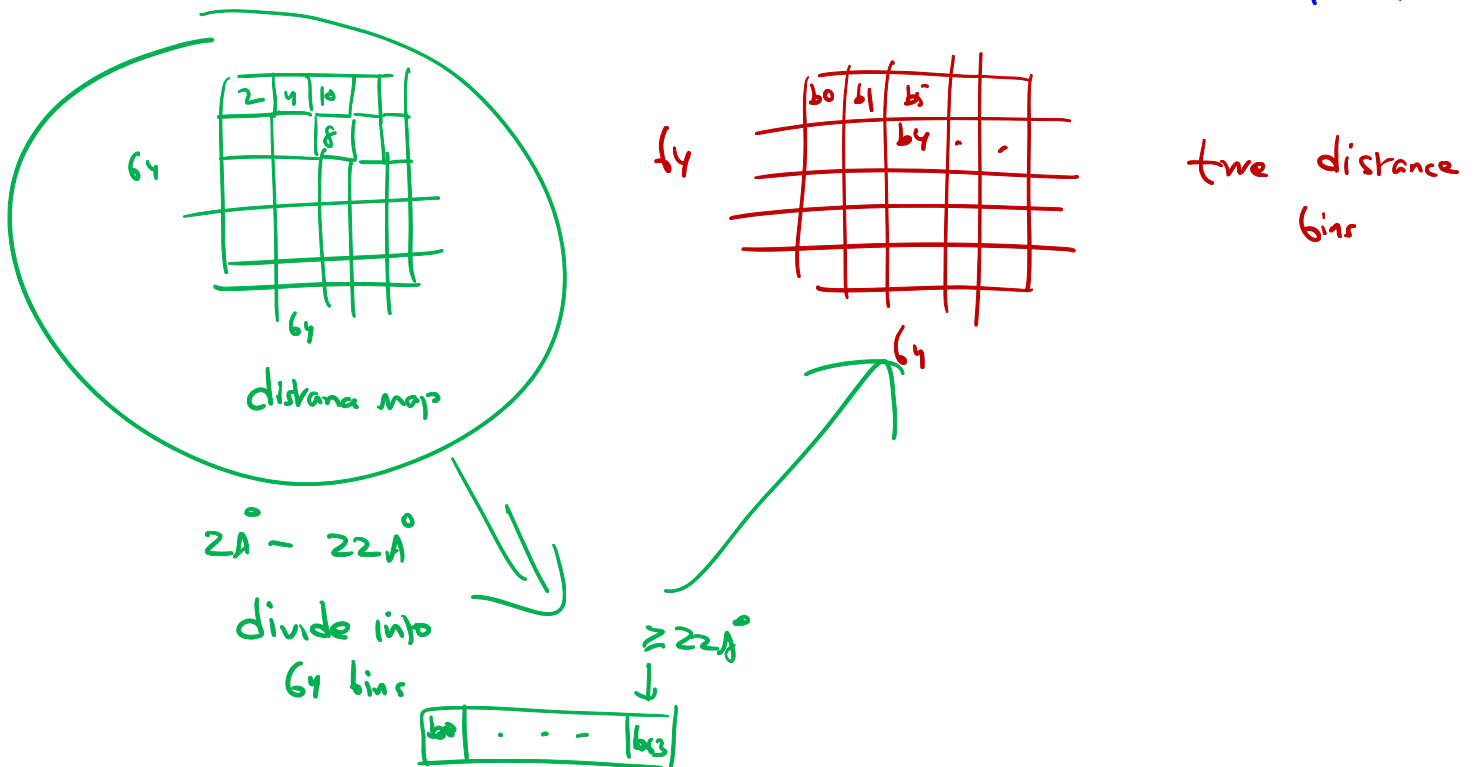
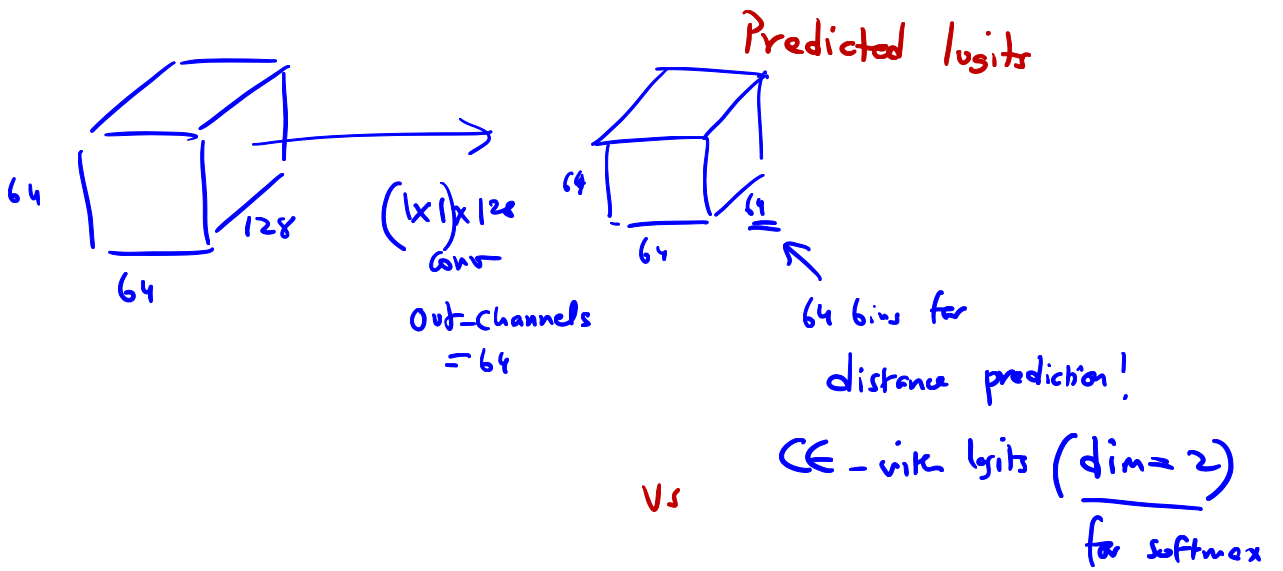


Dilatate G_{UV} ($d \equiv$ dialation)

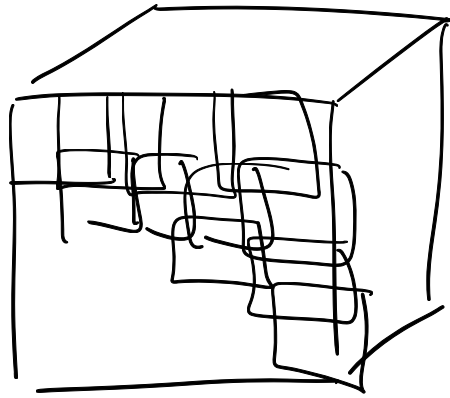




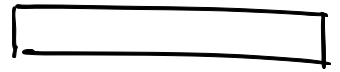
220 dilation blocks!



testing



$i, j =$



multiple
predictions

$$\max_k \left\{ P(b_k | i, j) \right\}$$

$k = 0 \dots 63$