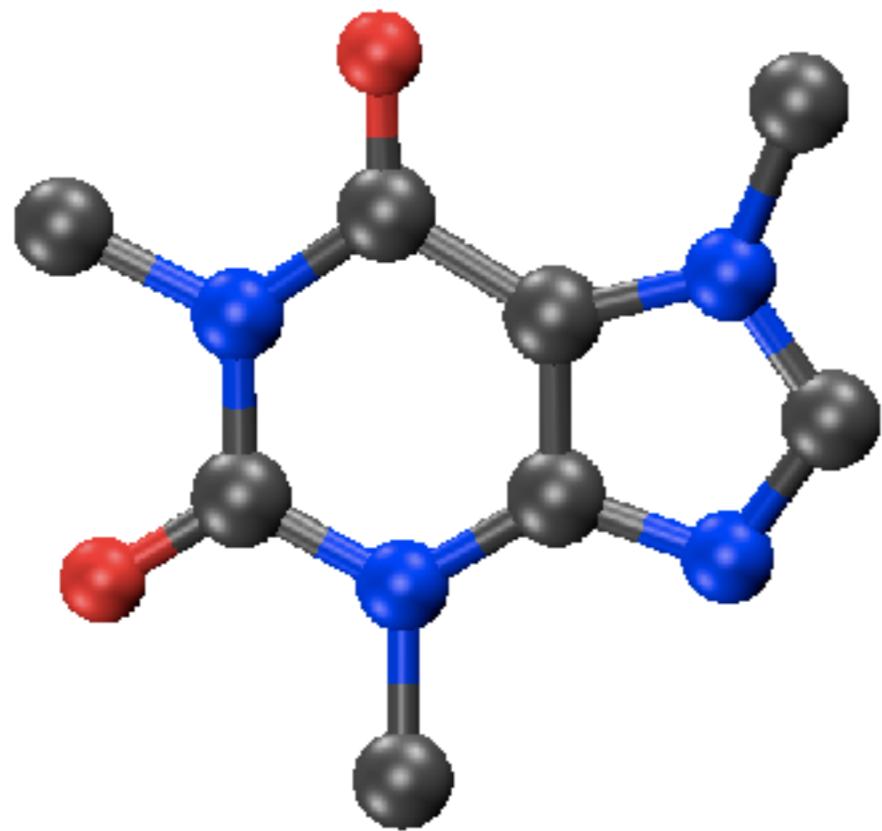
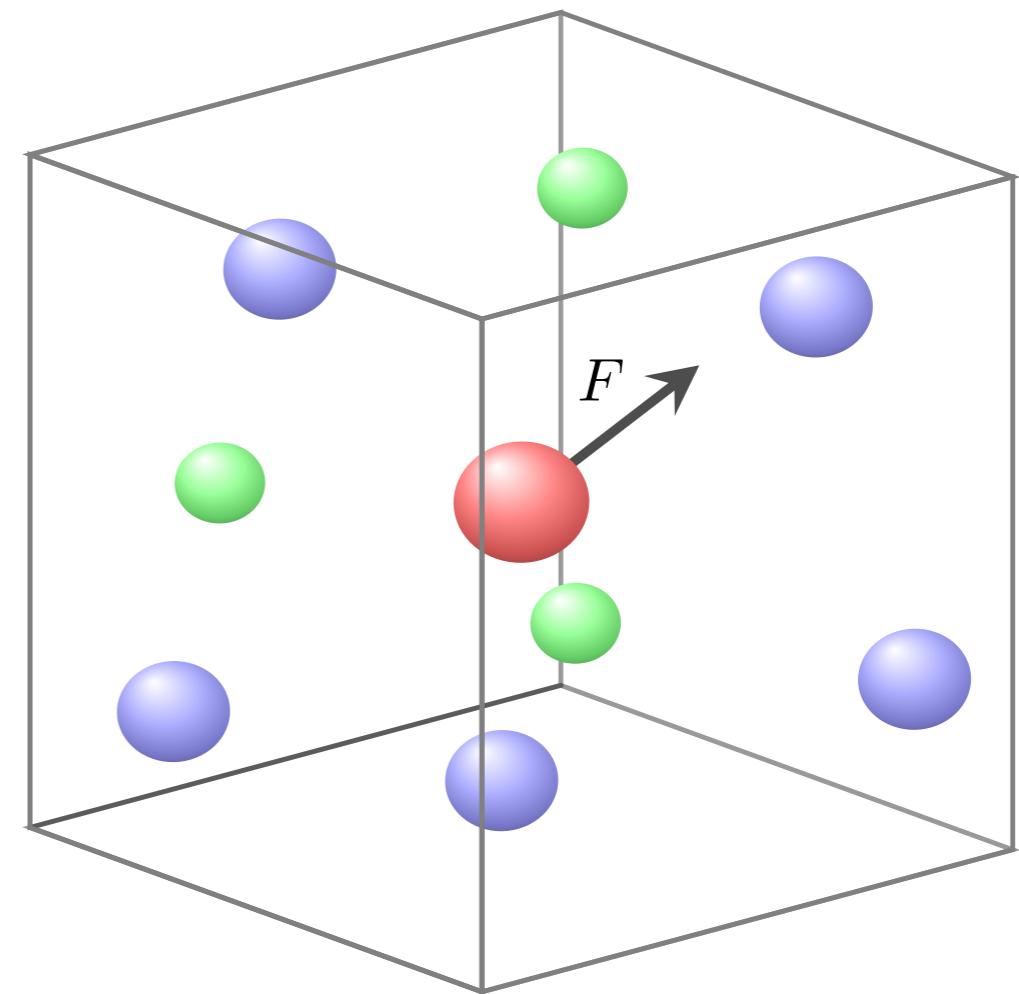


# Covariant neural networks for learning physical systems

Risi Kondor  
The University of Chicago



$$\phi(G)$$



$$F(\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_m)$$

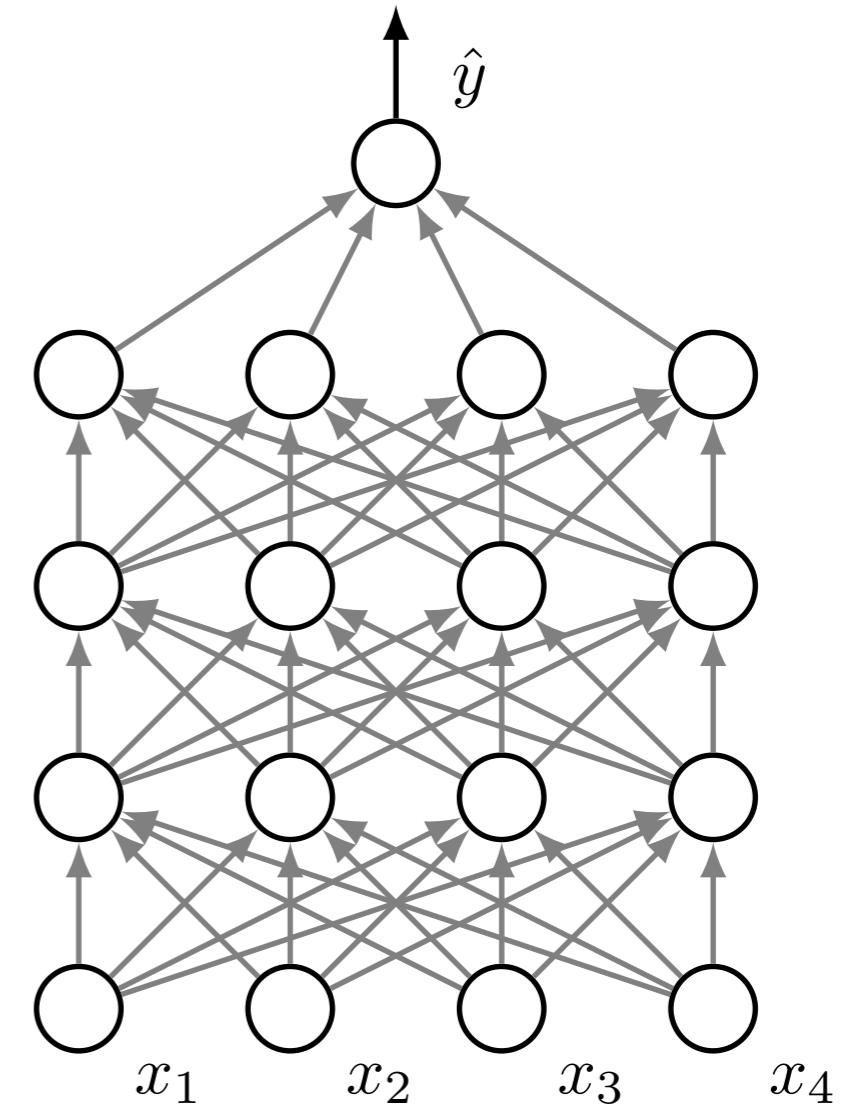
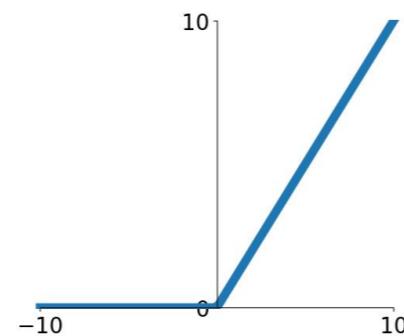
# 1. Multiscale structure and equivariance

# Feed-forward Neural Networks

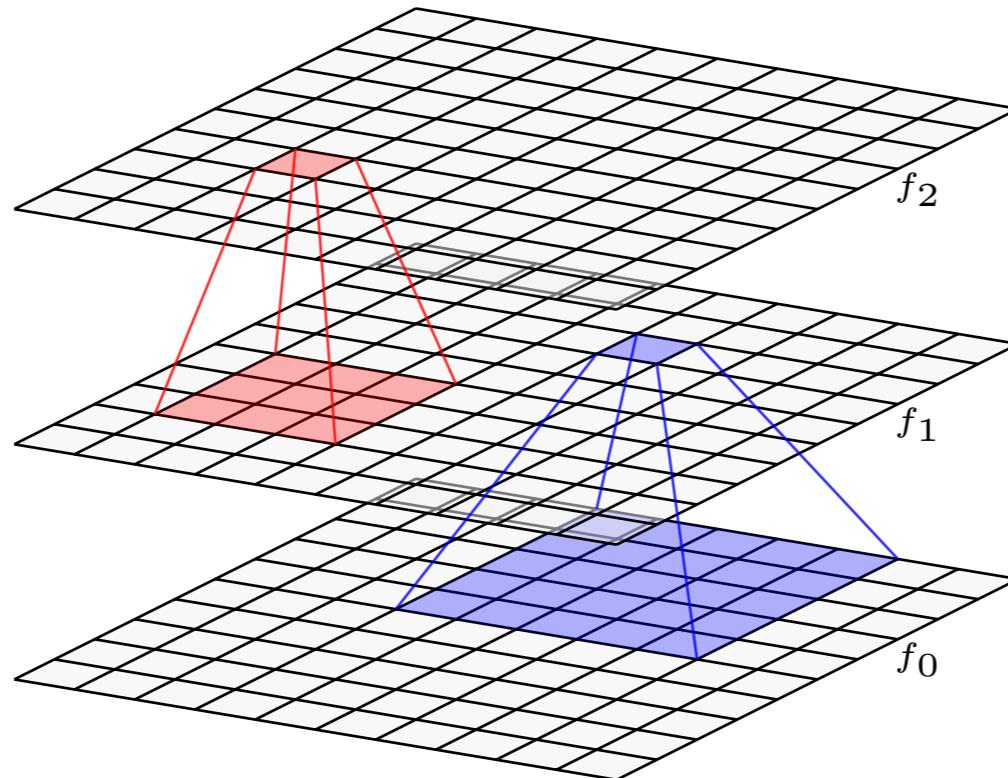
$$f_{\text{out}} = \sigma \left( \sum_i w_i f_{\text{in}}^{(i)} + b \right)$$

Common choice of nonlinearity:

$$\sigma_{\text{ReLU}}(z) = \max(0, x)$$



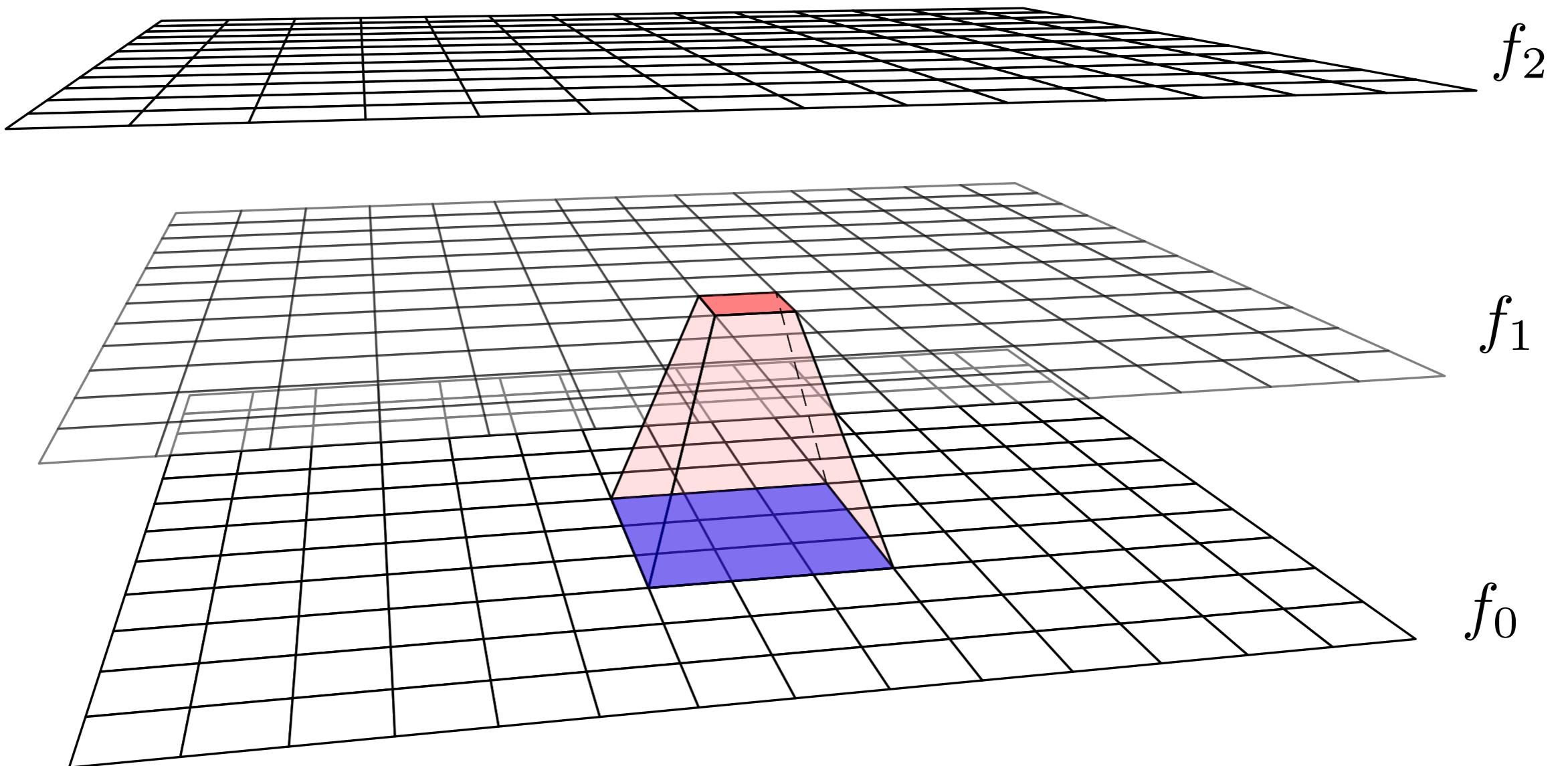
# Convolutional Neural Networks



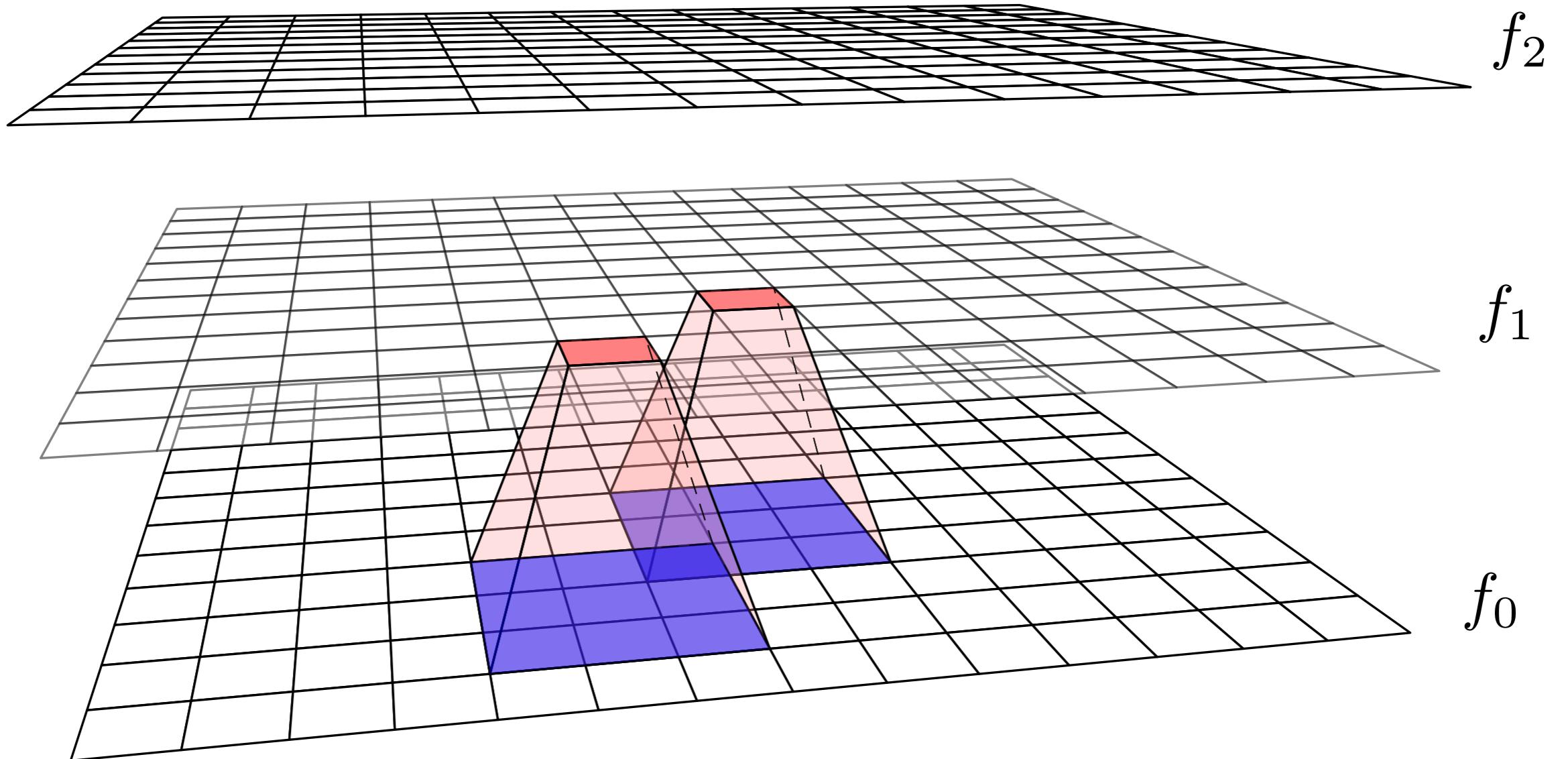
$$\phi_\ell(f_{\ell-1}) = (f_{\ell-1} * g_\ell)(x) = \sum_{y \in \mathbb{Z}^2} f_{\ell-1}(x - y) g_\ell(y)$$

Filter at layer  $\ell$

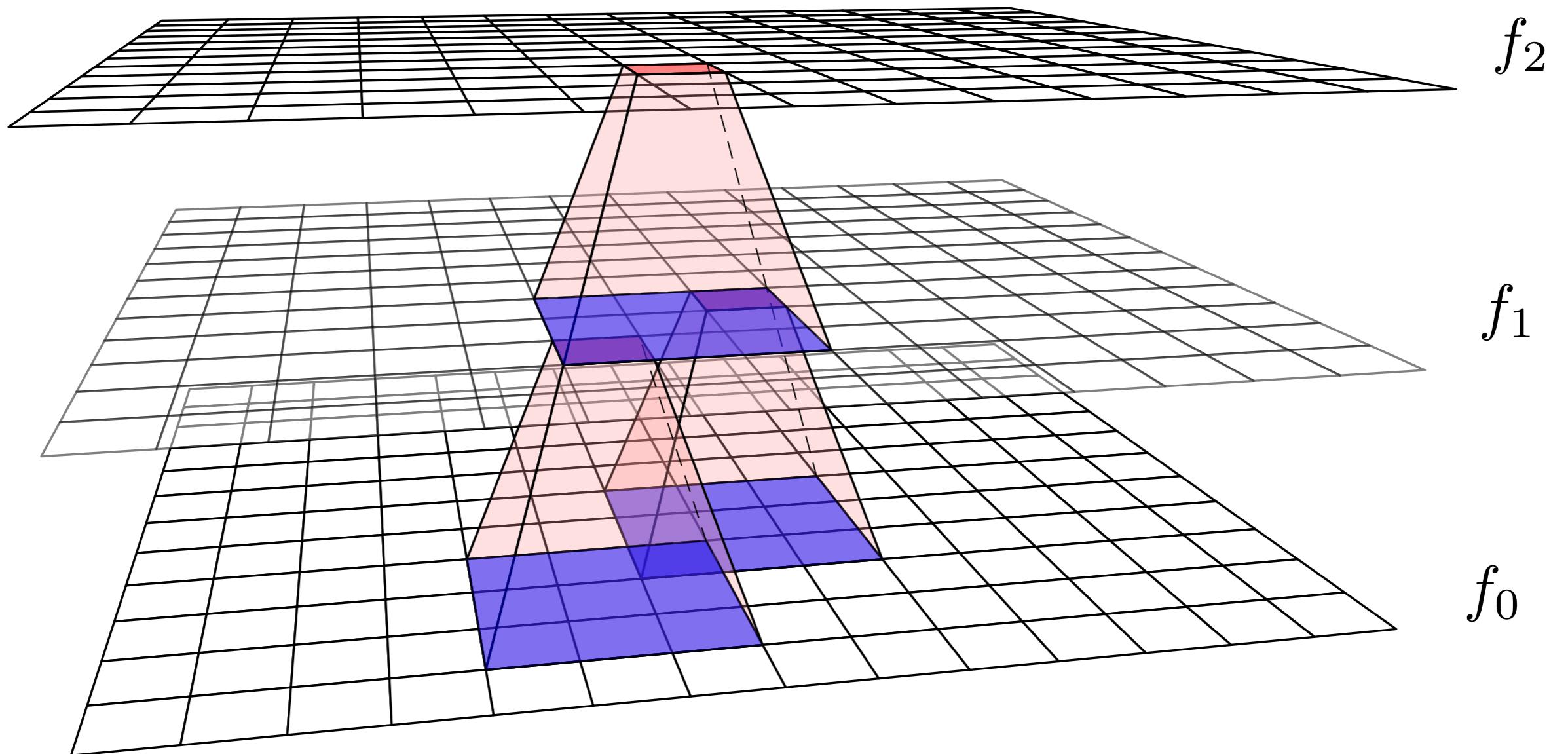
# Multiscale structure



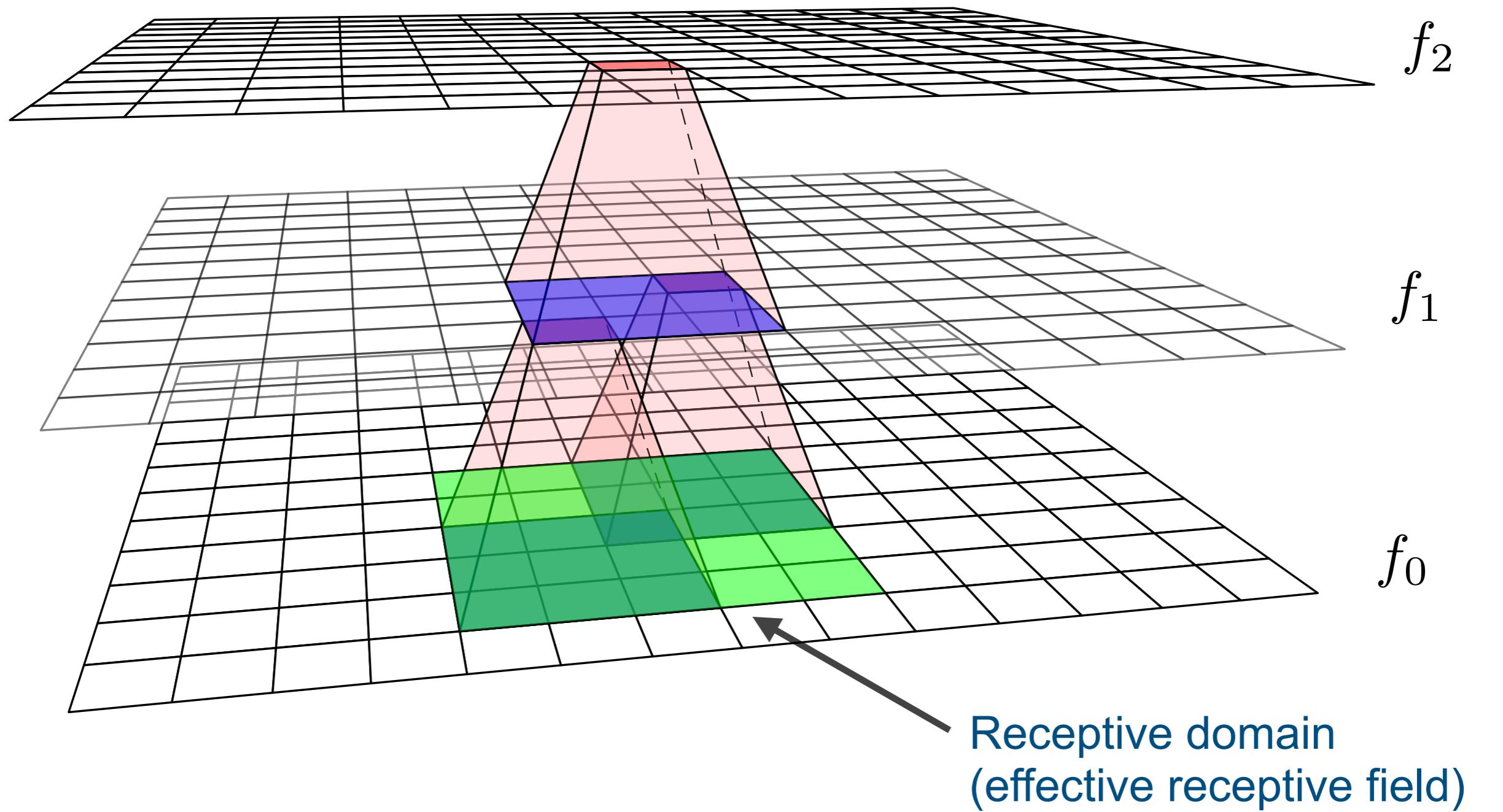
# Multiscale structure



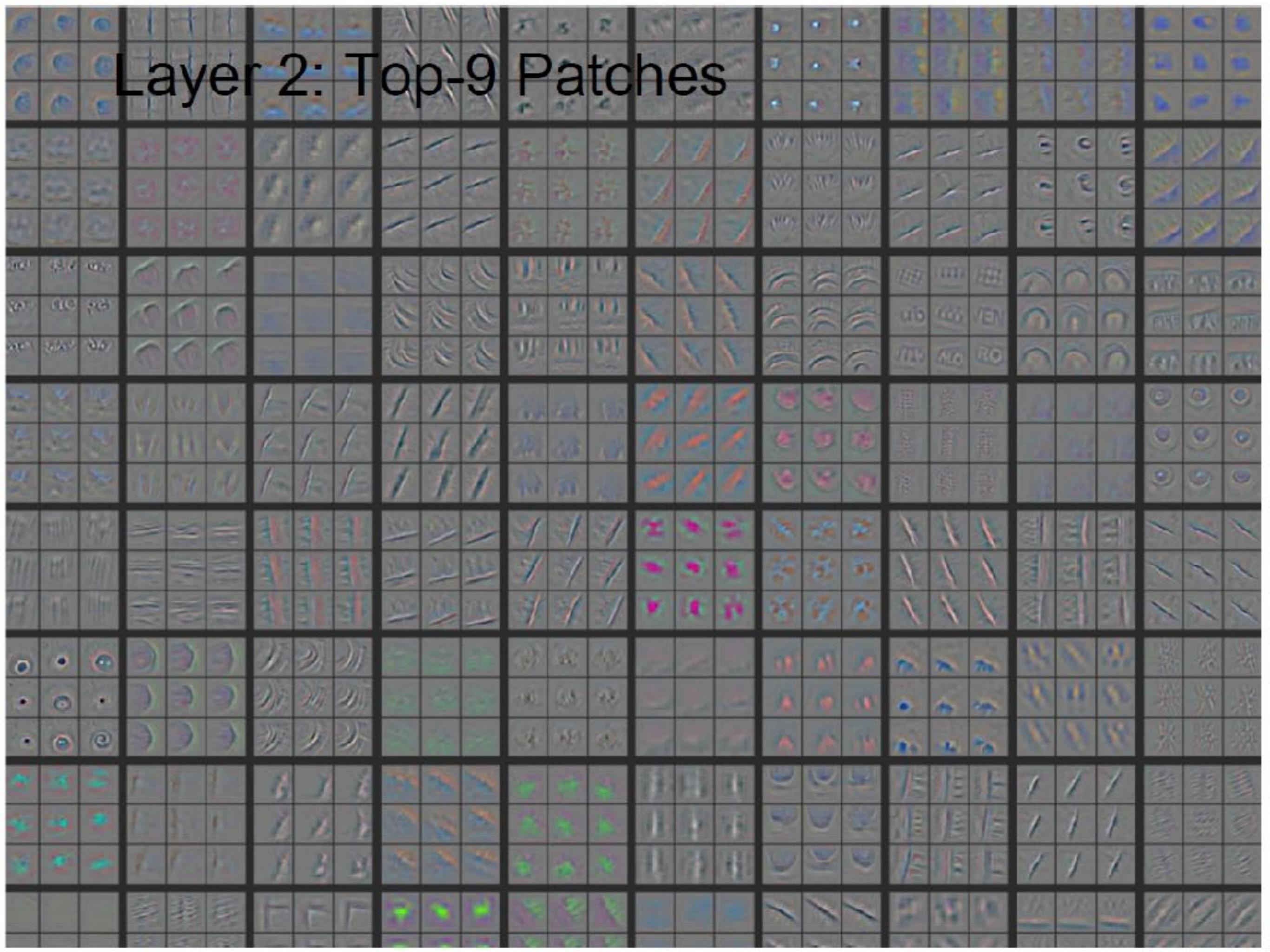
# Multiscale structure



# Multiscale structure



## Layer 2: Top-9 Patches



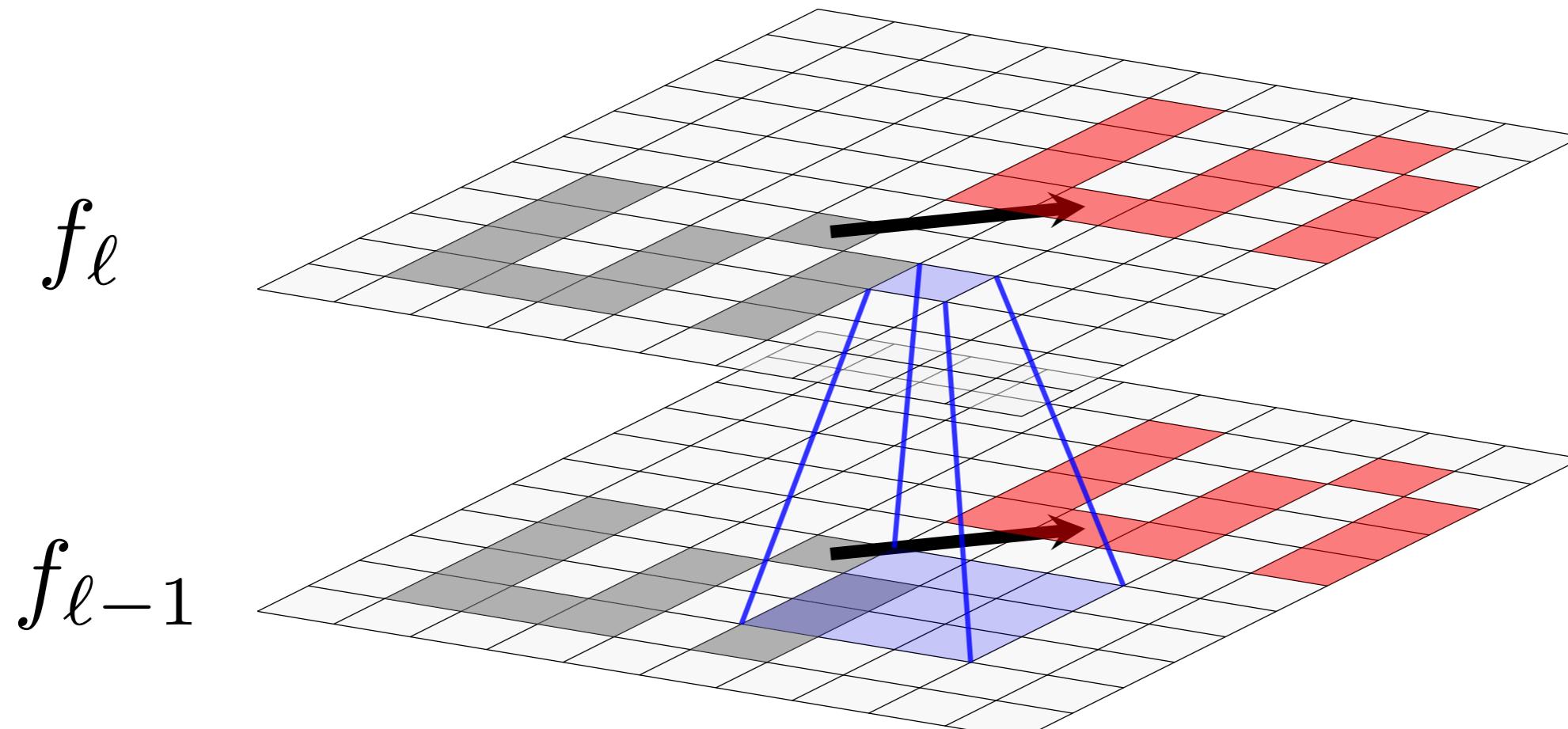
# Layer 3: Top-9 Patches



# Layer 4: Top-9 Patches

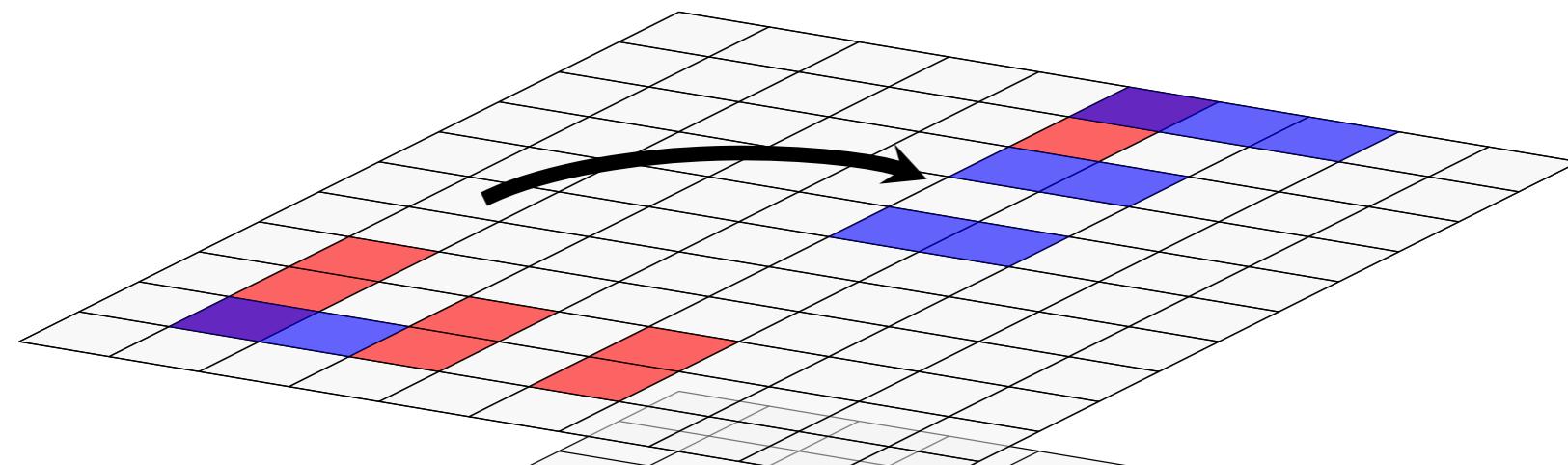


# Equivariance to translations

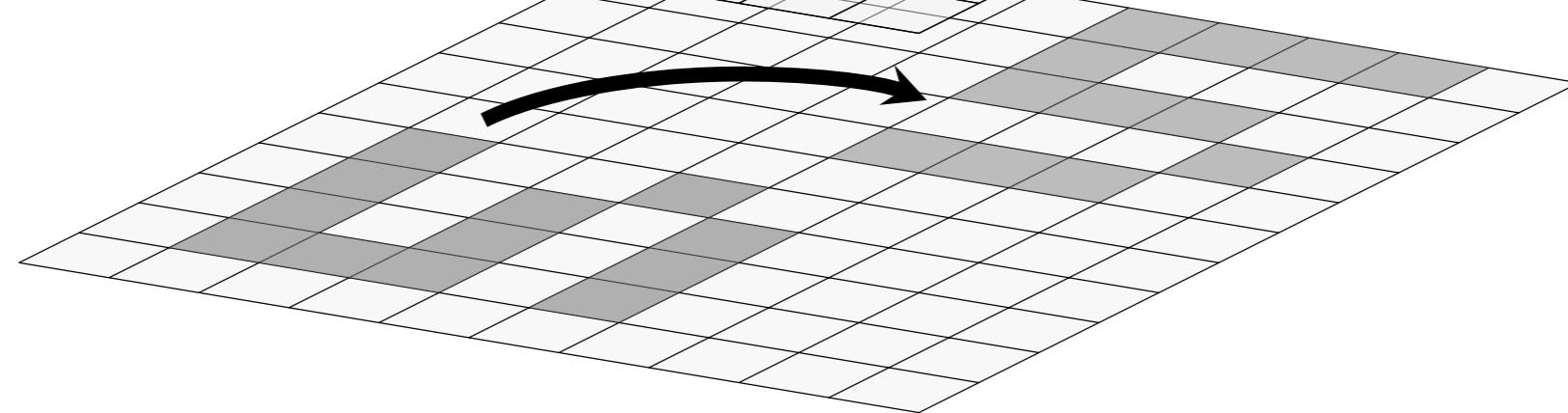


# Steerability

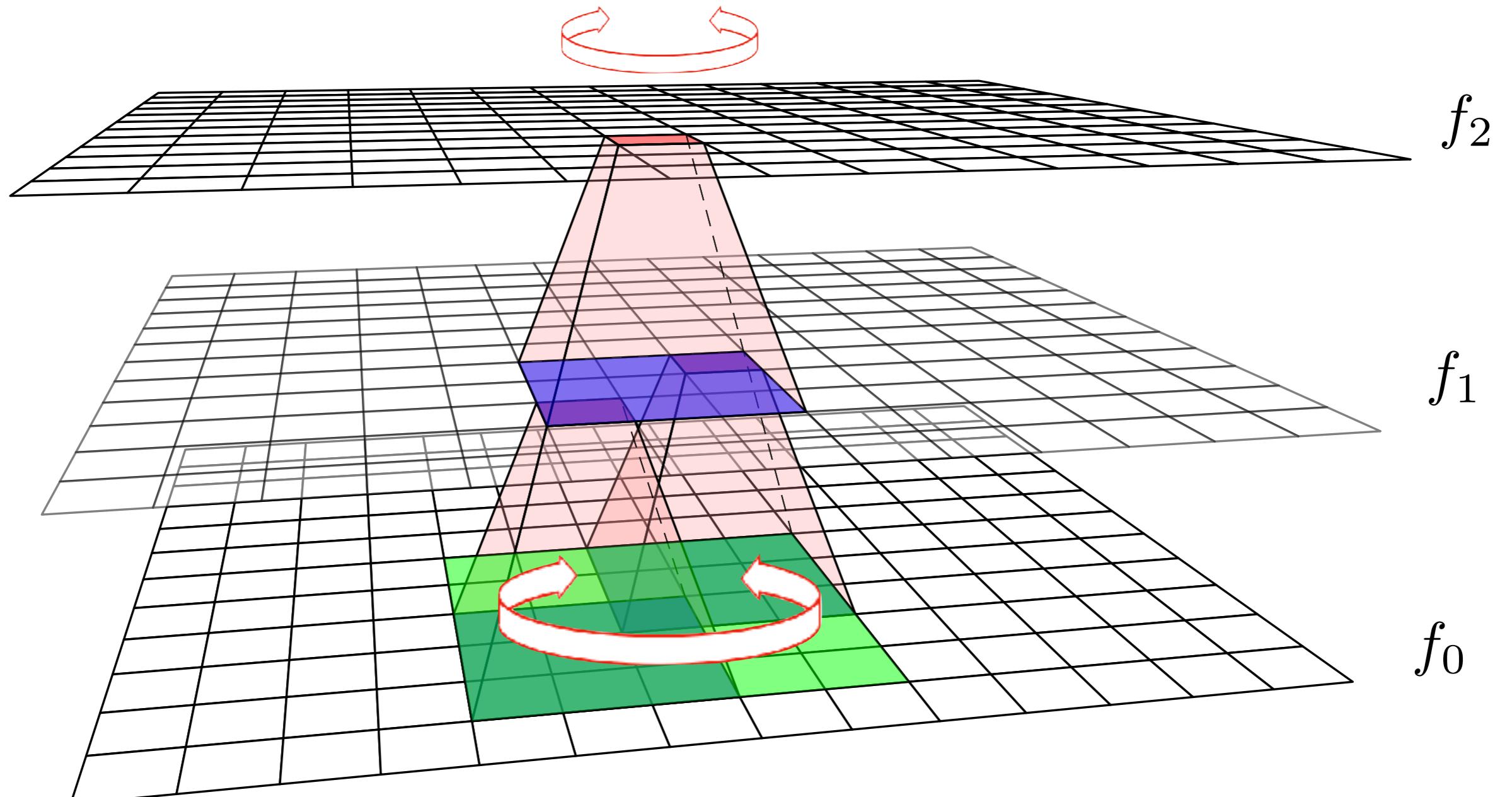
$f_\ell$



$f_{\ell-1}$

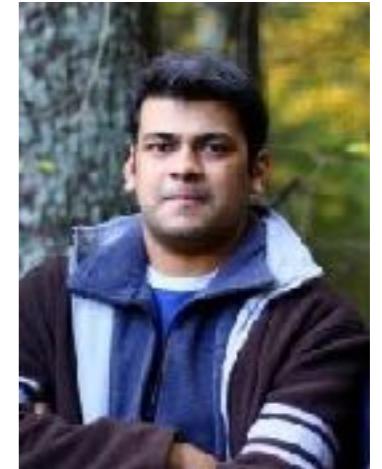


# Equivariance (covariance)



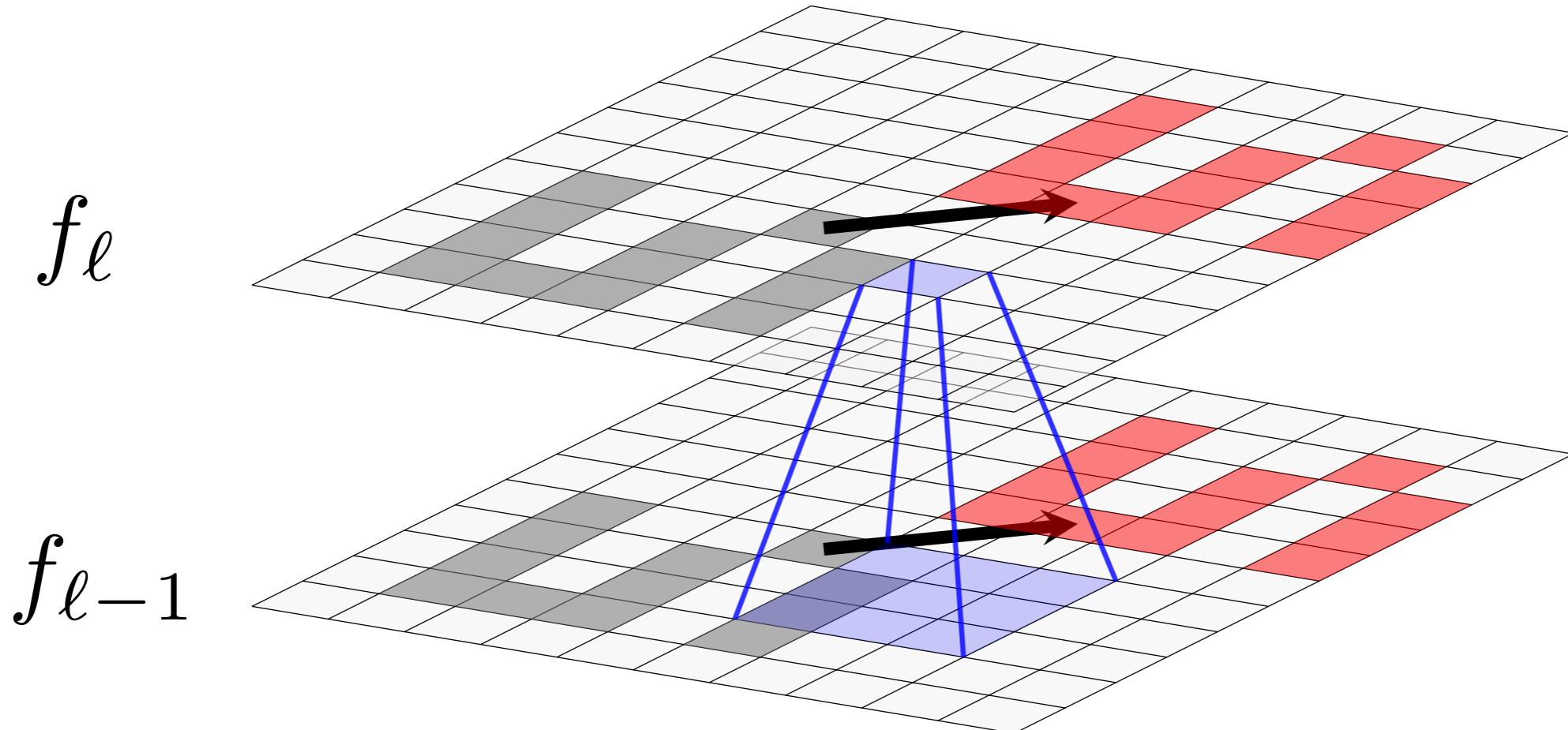
[Cohen & Welling, 2016]

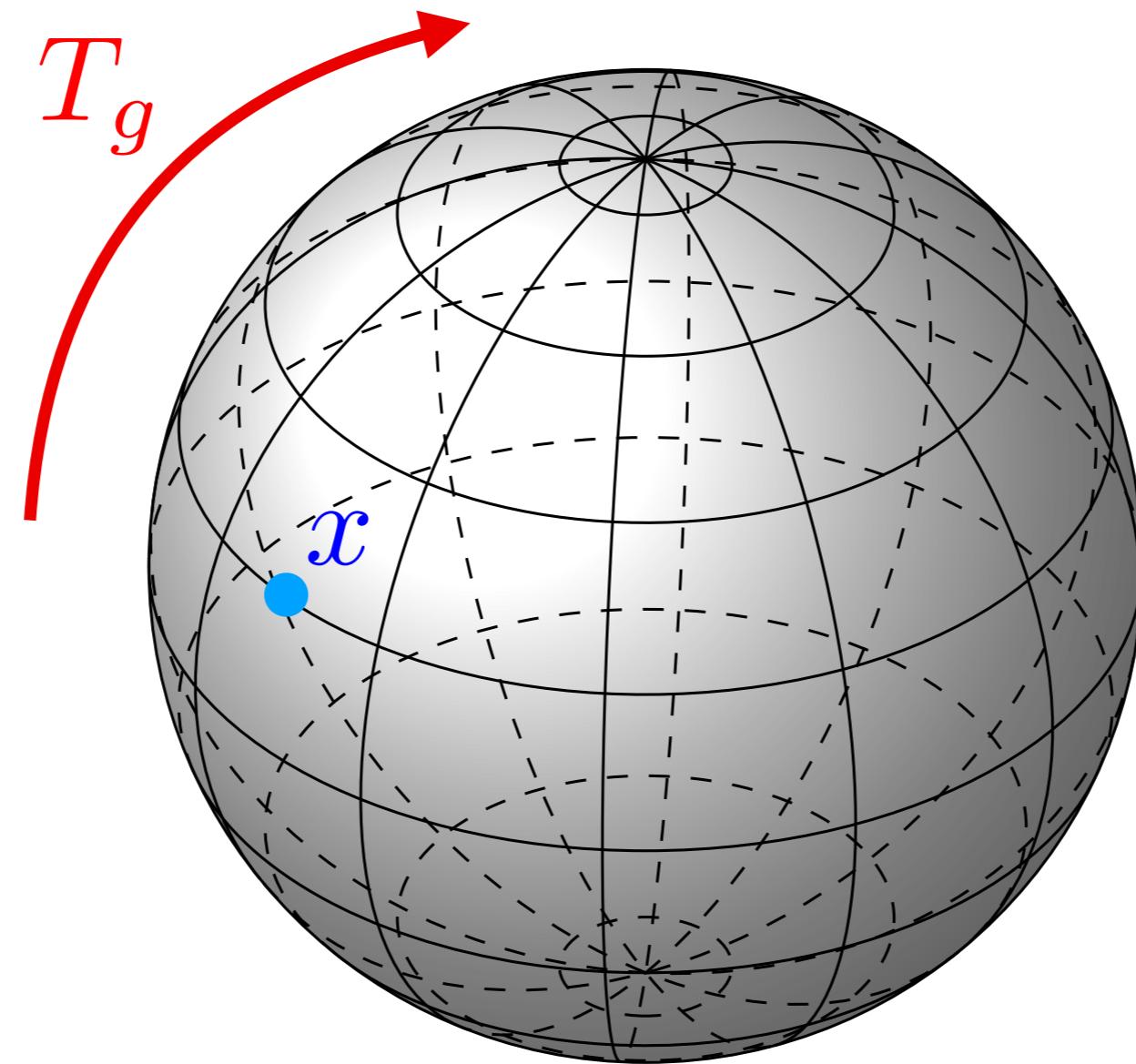
What is the analog of convolutions on graphs?



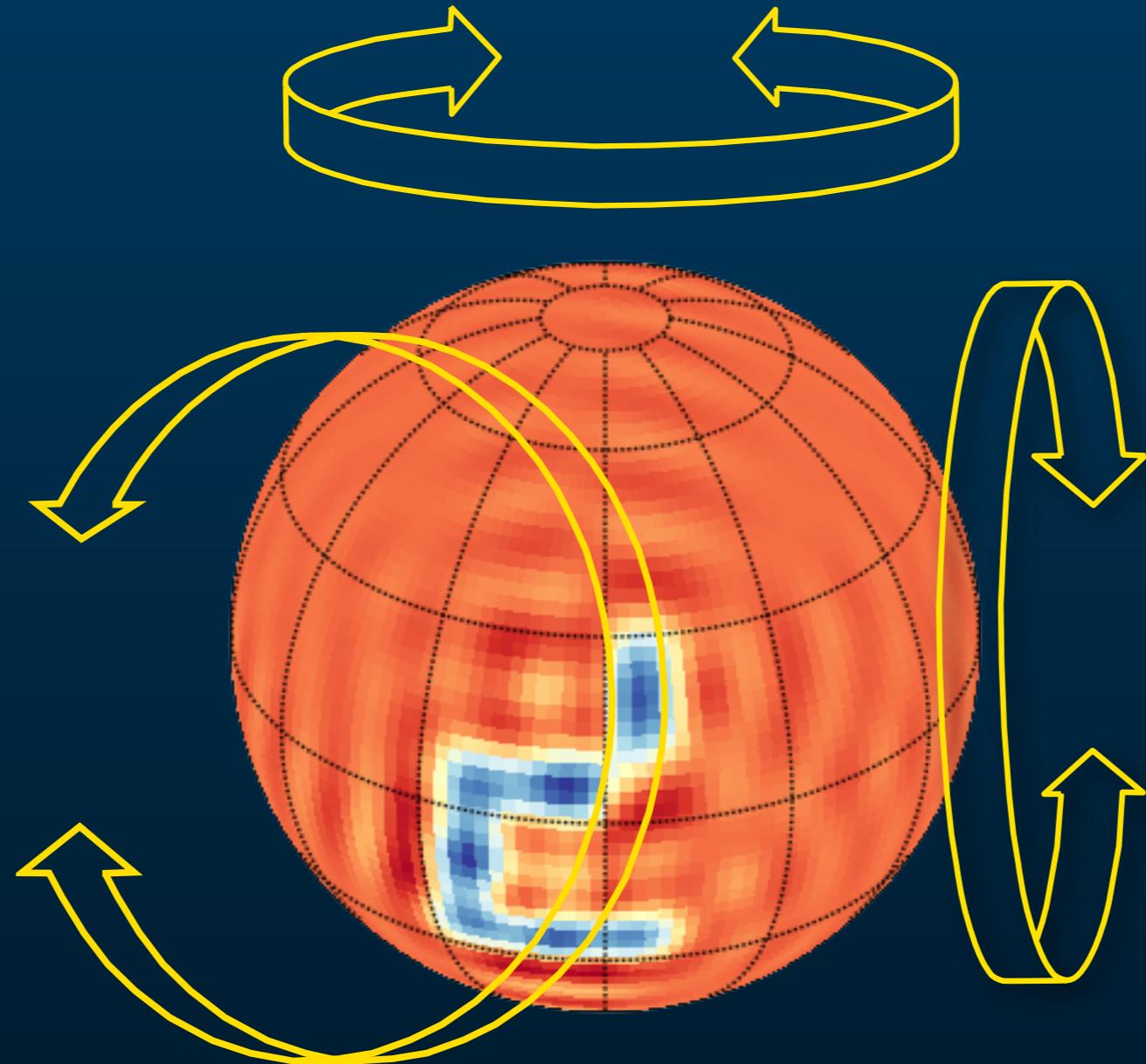
## 2. Theory

# How do we generalize convolution to the action of general (compact) groups?

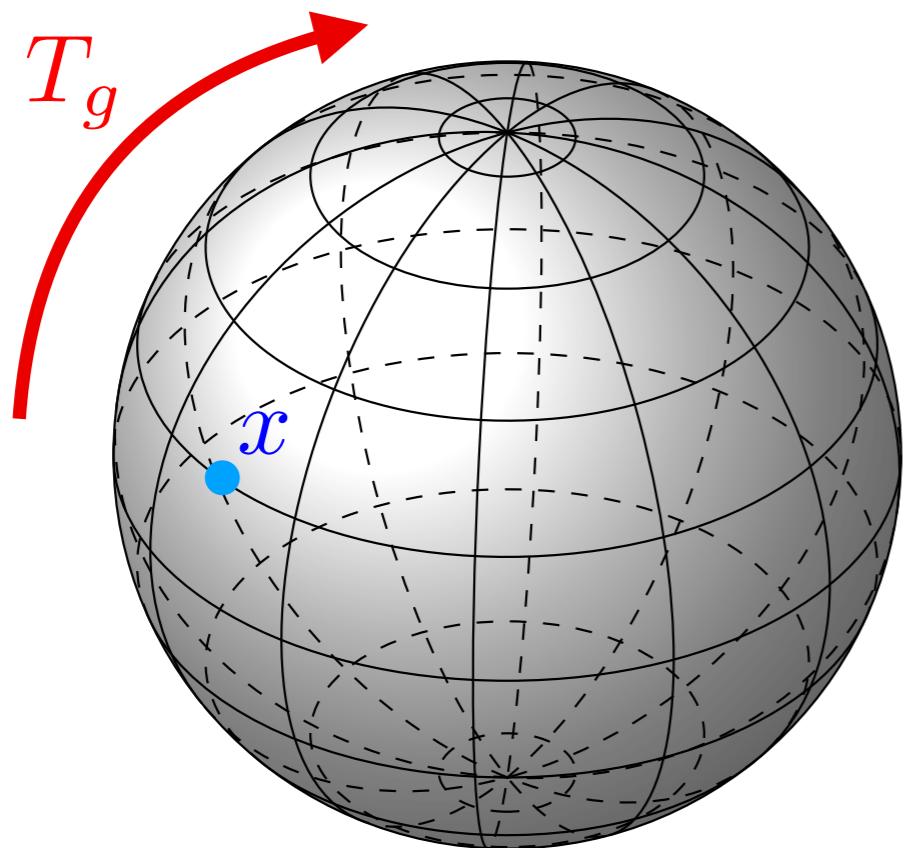




[Cohen, Geiger, Köhler & Welling, 2018]

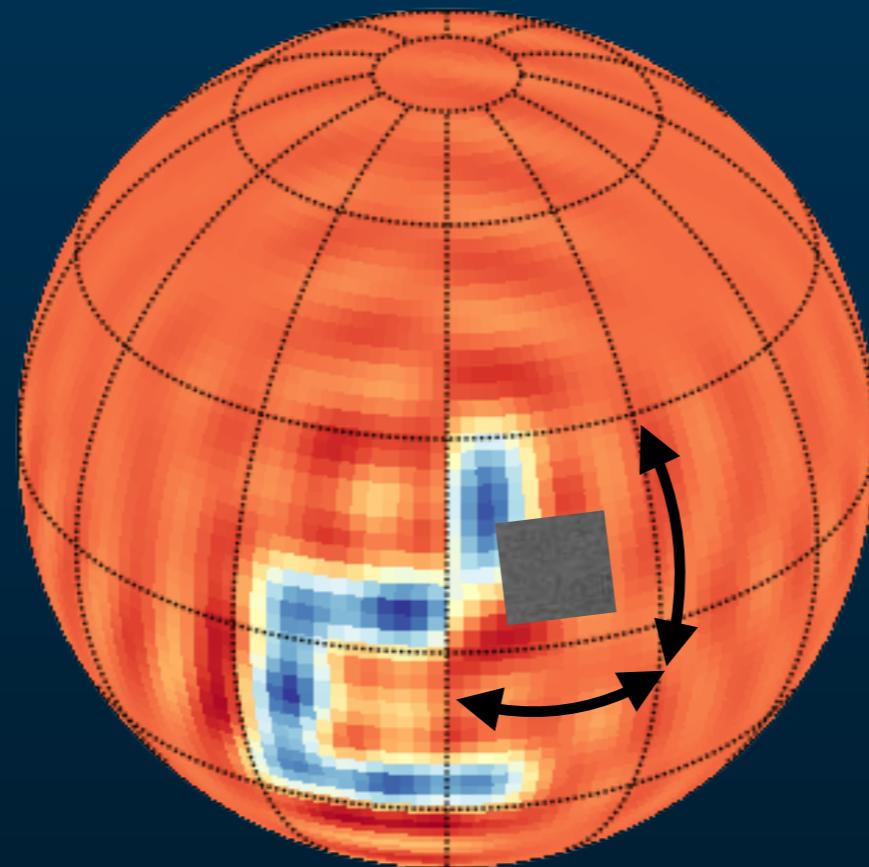


# Group actions



1. Our function lives on a space  $\mathcal{X}$   
$$f: \mathcal{X} \rightarrow \mathbb{C}$$
2. We have a group  $G$  acting on  $\mathcal{X}$   
$$x \mapsto T_g(x)$$
3. This induces an action on functions

$$f \xrightarrow{T_g} f' \quad f'(x) = f(T_g^{-1}(x))$$



$$(h \star f)(R) = \frac{1}{4\pi} \int_0^{2\pi} \int_{-\pi}^{\pi} [h_R(\theta, \phi)]^* f(\theta, \phi) \cos \theta d\theta d\phi$$

# Equivariance



$$\begin{array}{ccc} L(\mathcal{X}_1) & \xrightarrow{\mathbb{T}_g^{(1)}} & L(\mathcal{X}_1) \\ \downarrow \phi & & \downarrow \phi \\ L(\mathcal{X}_2) & \xrightarrow{\mathbb{T}_g^{(2)}} & L(\mathcal{X}_2) \end{array}$$

$$(f*g)(u)=\int_G f(uv^{-1})\,g(v)\,d\mu(v)$$

$$(f*g)(u)=\int_G f{\restriction} G(uv^{-1})\,g{\restriction} G(v)\,d\mu(v)$$

# Theorem

A feed-forward neural network is equivariant to the action of a compact group  $G$  if and only if the linear operation in each layer is of the form

$$\phi_\ell(f_{\ell-1}) = f_{\ell-1} * g_\ell.$$

where  $*$  denotes the generalization of convolution to homogeneous space of compact groups, defined

$$(f * g)(u) = \int_G f \uparrow^G(uv^{-1}) g \uparrow^G(v) d\mu(v)$$

# Convolution on groups



$$\widehat{f}(k) = \int f(x) e^{-ikx} dx$$

$$\widehat{(f * \chi)}(k) = \widehat{f}(k) \cdot \widehat{\chi}(k)$$

$$\widehat{f}(\rho) = \int f(x) \rho(x) d\mu(x)$$

$$(f * g)(u) = \int_G f(uv^{-1}) g(v) d\mu(v)$$

$$\widehat{(f * \chi)}(\rho) = \widehat{f}(\rho) \cdot \widehat{\chi}(\rho)$$



$$\rho_0 \qquad \rho_1 \qquad \rho_2 \qquad \rho_p$$

A diagram illustrating the decomposition of a space  $L(\mathcal{X})$  into a direct sum of subspaces  $V_i$ . The top row contains labels  $\rho_0, \rho_1, \rho_2, \dots, \rho_p$ . Below this, the equation  $L(\mathcal{X}) = V_0 \oplus V_1 \oplus V_2 \oplus \dots \oplus V_p$  is written. Four arrows originate from the labels  $\rho_0, \rho_1, \rho_2, \rho_p$  and point downwards towards the corresponding subspaces  $V_0, V_1, V_2, V_p$  in the equation below.

$$L(\mathcal{X}) = V_0 \oplus V_1 \oplus V_2 \oplus \dots \oplus V_p$$

$$V_i=W_i^1\oplus W_i^2\oplus\ldots\oplus W_i^{m_i}$$

$$\widehat{f}(\rho_i) = \int_G f(u) \rho_i(u) d\mu(u) \quad i = 0, 1, 2, \dots$$

$$\widehat{f * g}(\rho_i) = \widehat{f}(\rho_i) \cdot \widehat{g}(\rho_i)$$



matrix multiplication

**Case 1:**  $f_{\ell-1}: G/H \rightarrow \mathbb{C}$   $f_\ell: G \rightarrow \mathbb{C}$

$$\left( \begin{array}{c} \text{[Gray rectangle]} \end{array} \right) = \left( \begin{array}{c} \text{[White rectangle with vertical gray bars]} \end{array} \right) \times \left( \begin{array}{c} \text{[White rectangle with horizontal gray bars]} \end{array} \right)$$

$$\widehat{f * g}(\rho)$$

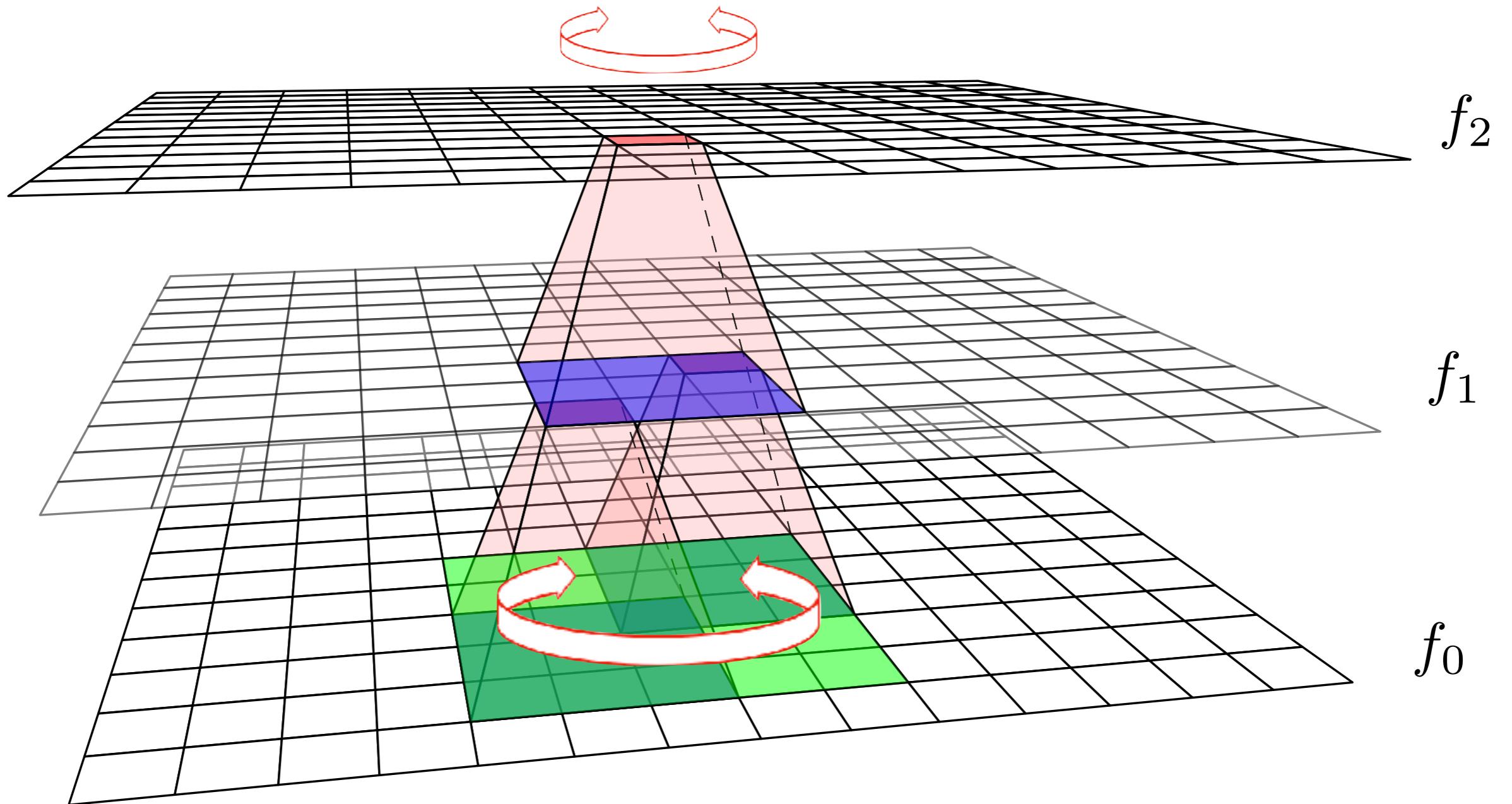
$$\widehat{f \uparrow G}(\rho)$$

$$\widehat{g \uparrow G}(\rho)$$

**Case 2:**  $f_{\ell-1}: G/H \rightarrow \mathbb{C}$        $f_\ell: G/K \rightarrow \mathbb{C}$

$$\left( \begin{array}{c|c|c|c|c} \hline & & & & \\ \hline \end{array} \right) = \left( \begin{array}{c|c|c|c} \hline & & & \\ \hline \end{array} \right) \times \left( \begin{array}{ccc} \textcolor{gray}{\square} & \textcolor{gray}{\square} & \textcolor{gray}{\square} \\ \vdots & \vdots & \vdots \\ \textcolor{gray}{\square} & \textcolor{gray}{\square} & \textcolor{gray}{\square} \end{array} \right)$$
$$\widehat{f * g}(\rho) \qquad \widehat{f \uparrow G}(\rho) \qquad \widehat{g \uparrow G}(\rho)$$

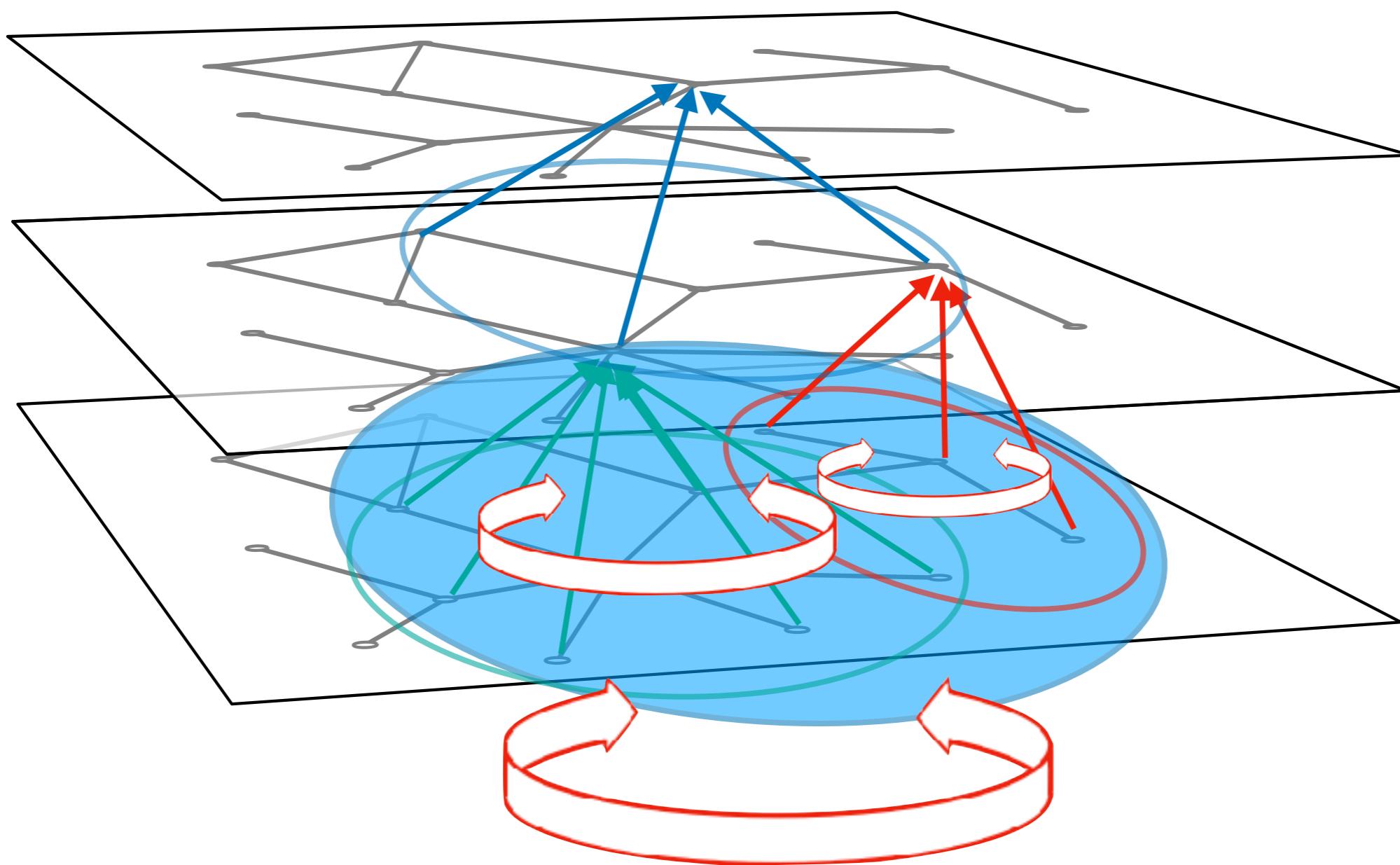
# Fourier Space Neural Networks



### 3. Covariant graph neural networks

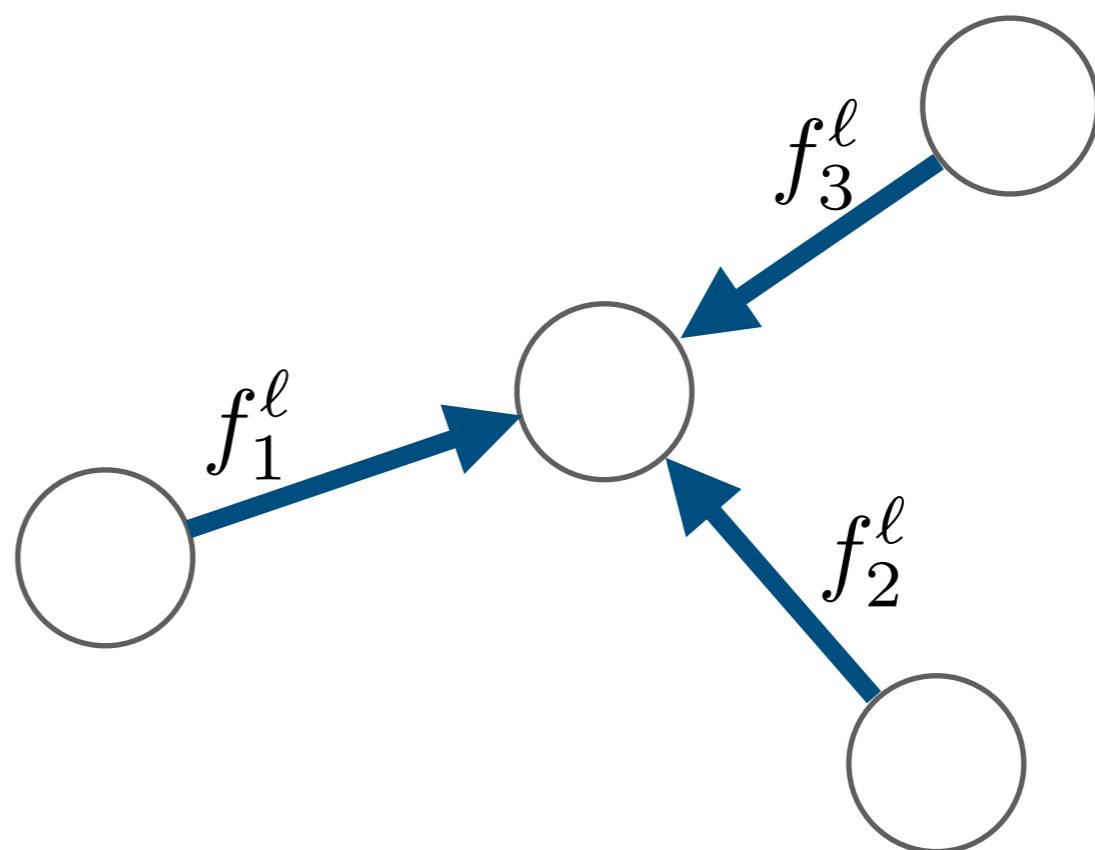
[K., Pan, Hy-Truong, Trivedi & Anderson]

# Multiscale structure in graphs



# Message passing neural networks (MPNNs)

$$f_i^{\ell+1} = \xi \left( W \sum_{j \in \mathcal{N}(i)} f_j^\ell + b \right)$$



[Gilmer et al, '17] [Krieger, '16] [Niepert, '16] [Duvenaud et al., '15]  
[Dai, Dai & Song, '16]

# Permutation covariant activations

How does a given activation transform when its receptive field is subjected to a permutation  $\sigma$ ?

Scalar activation

(0th order covariant)

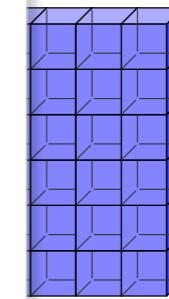
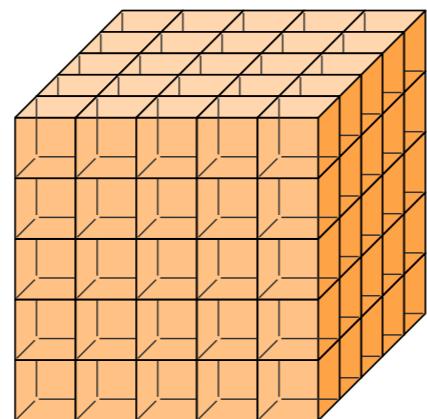
General tensor activation  
(kth order covariant)

Vector activation

(1st order covariant)

Matrix activation

(2nd order covariant)

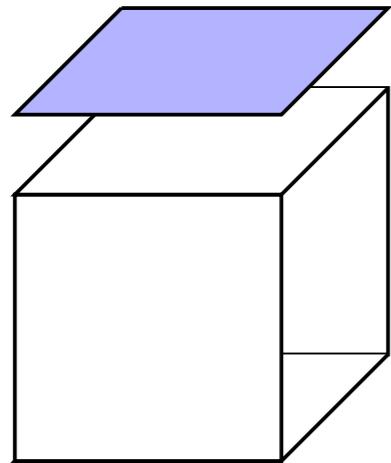


$$P_\sigma F_i P_\sigma^\top$$

$$F_{i_1, i_2, \dots, i_k} \xrightarrow{\sigma} [P_\sigma]_{i_1}^{j_1} [P_\sigma]_{i_2}^{j_2} \dots [P_\sigma]_{i_k}^{j_k} F_{j_1, j_2, \dots, j_k}$$

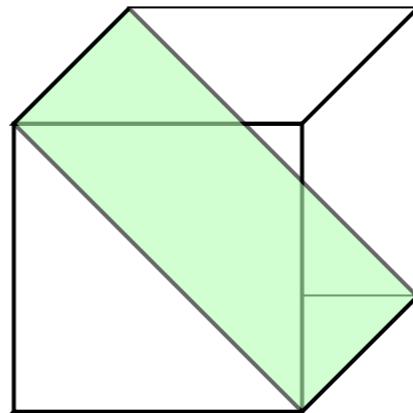
# Permutation covariant operations:

1. Projections



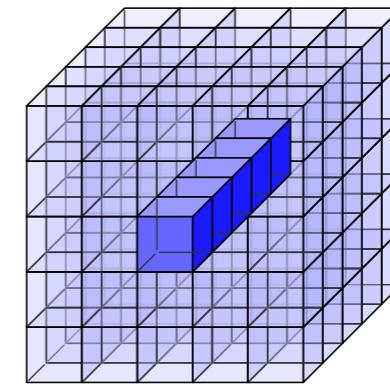
$$C_{i,j} = \sum_a A_{a,i,j}$$

2. Diagonals



$$C_{i,j} = \sum_i A_{i,i,j}$$

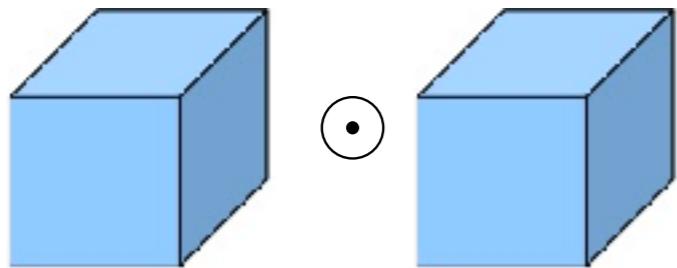
3. Contractions



$$C_k = \sum_{i,j} A_{i,j,k}$$

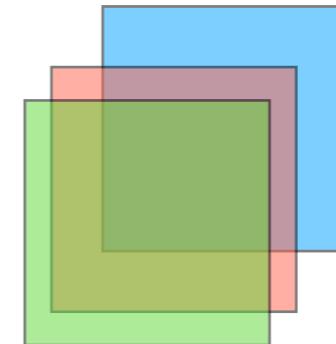
# Permutation covariant operations:

## 4. Hadamard products



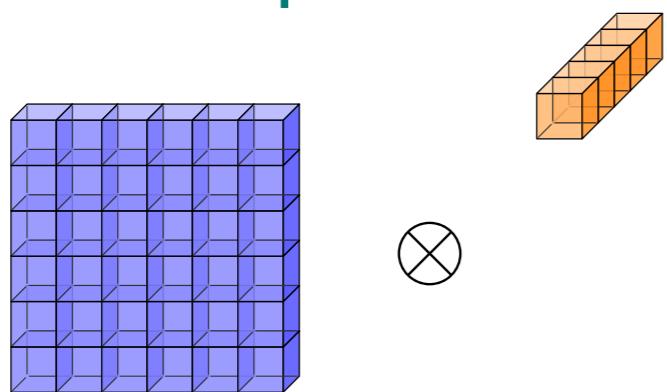
$$C_{i,j,k} = A_{i,j,k} B_{i,j,k}$$

## 5. Stacking



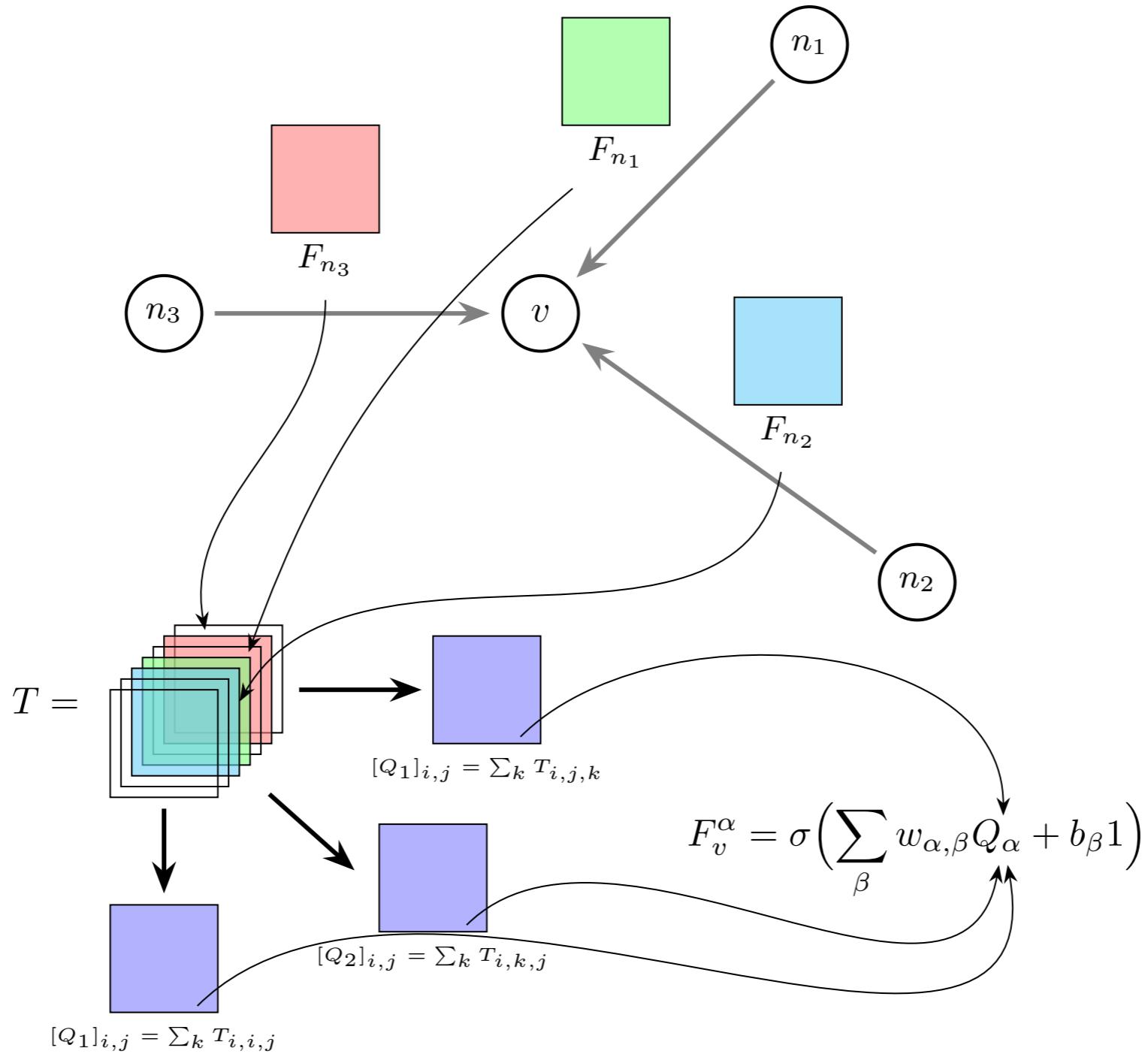
$$C_{i,j,k} = A_{i,j}^{(k)}$$

## 6. Tensor products



$$C_{i,j,k} = A_{i,j} B_k$$

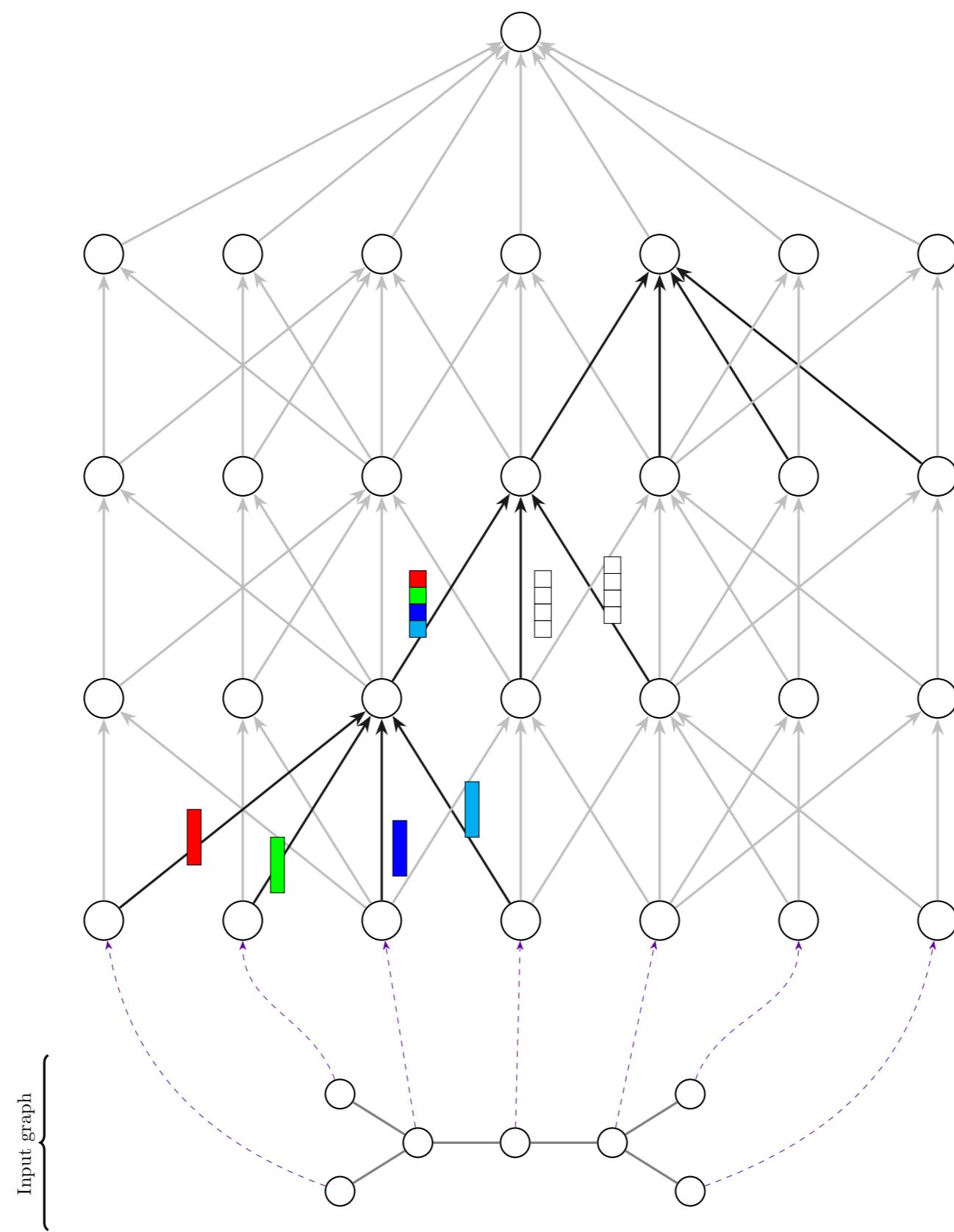
# Second order vertex aggregation



# Harvard Clean Energy Project results

Method	Train MAE	Train RMSE	Test MAE	Test RMSE
Lasso	0.863	1.190	0.867	1.437
Ridge Regression	0.849	1.164	0.854	1.376
Random Forest	0.999	1.331	1.004	1.799
Gradient Boosted Tree	0.676	0.939	0.704	1.005
Weisfeiler-Lehman Graph Kernel	0.805	1.111	0.805	1.096
Neural Graph Fingerprint	0.848	1.187	0.851	1.177
Learning Convolution Neural Network	0.704	0.972	0.718	0.973
CCN 2D	0.562	0.773	0.570	0.773

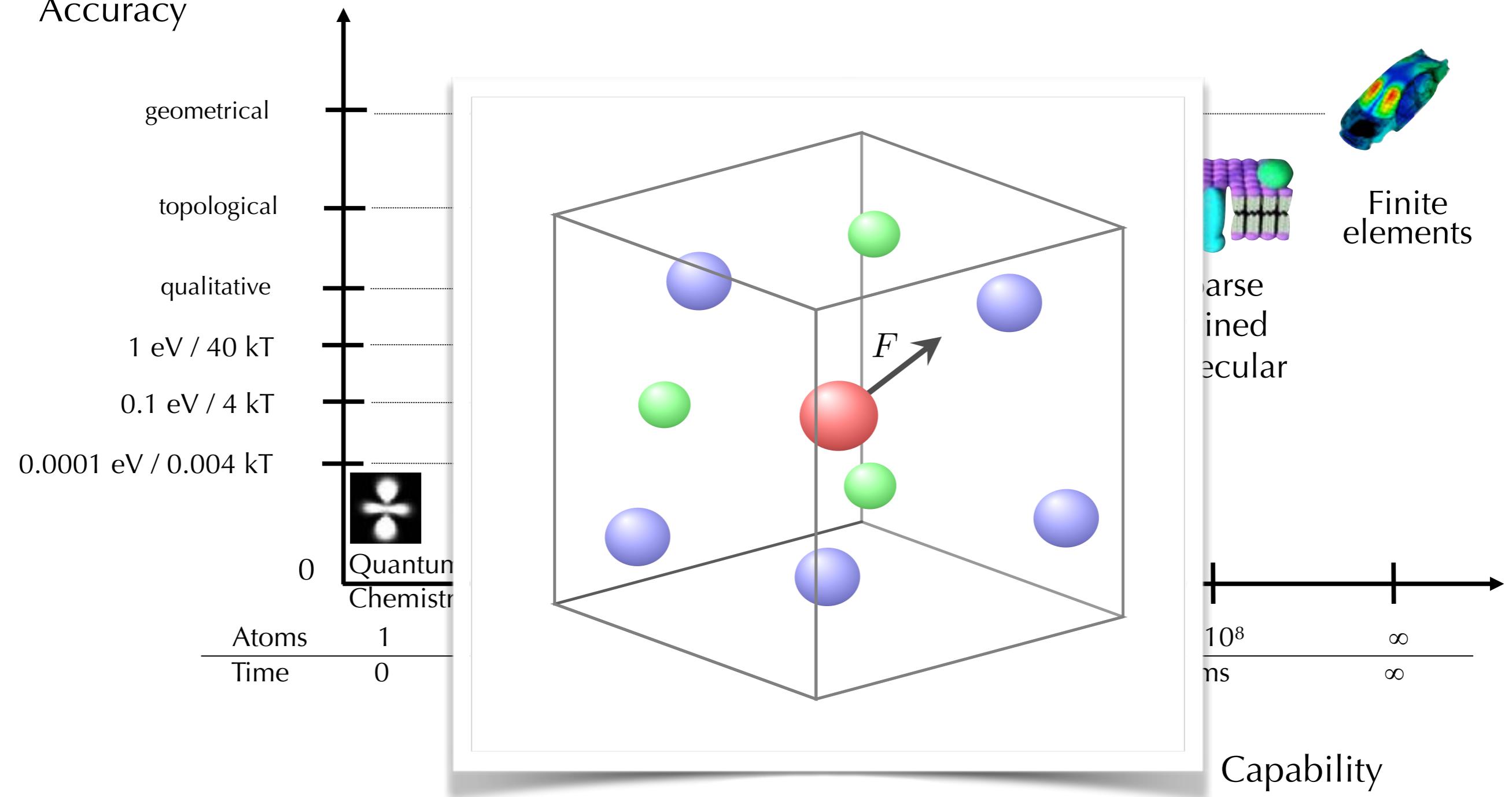
[Duvenaud et al., 2015] [Kriege, 2016] [Niepert, 2016]  
[Hachmann et al., 2011]



## 4. $\text{SO}(3)$ covariant neural networks

# Multiple scales of materials modelling

Accuracy



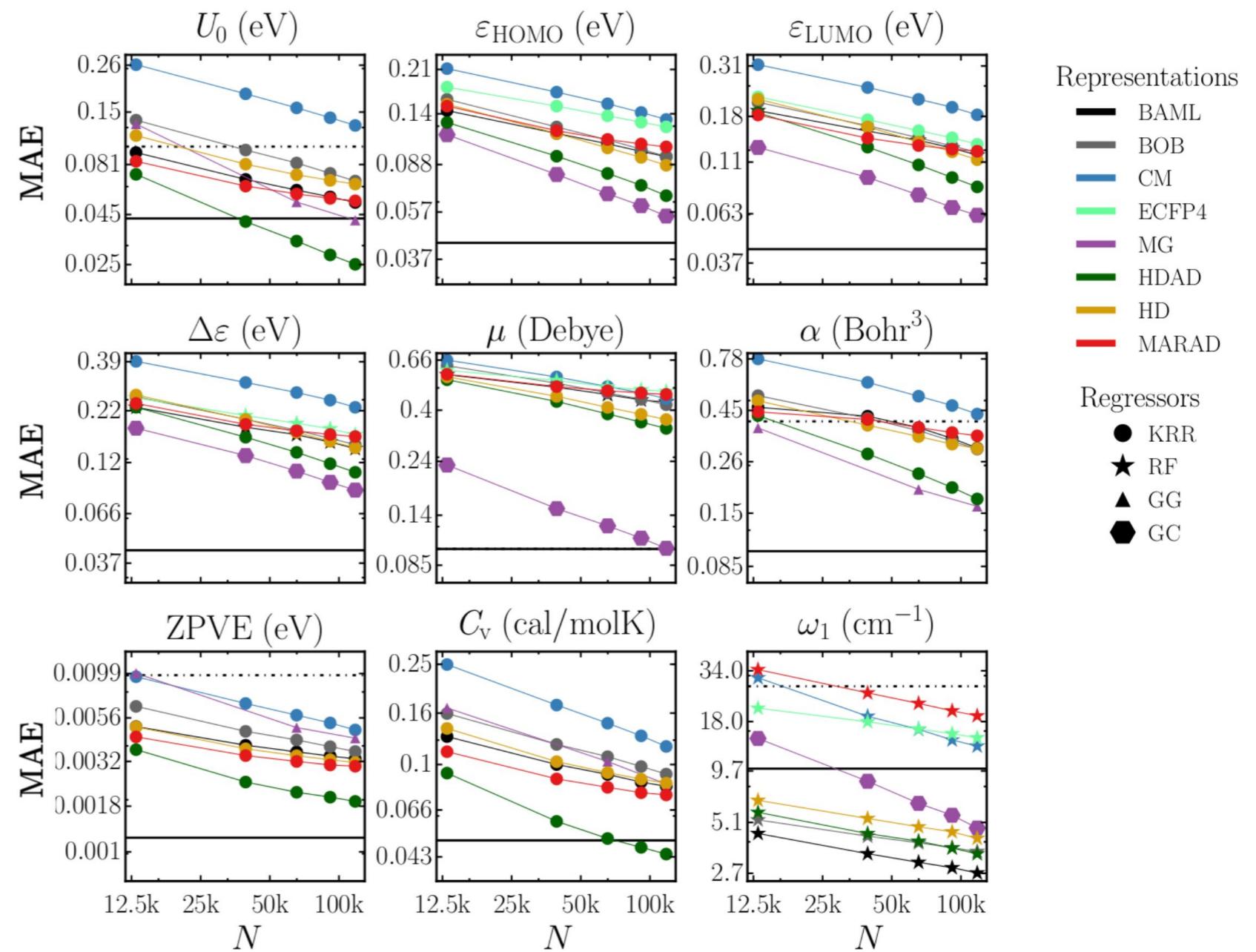
Capability

[Csányi]

# The GDB-9 dataset

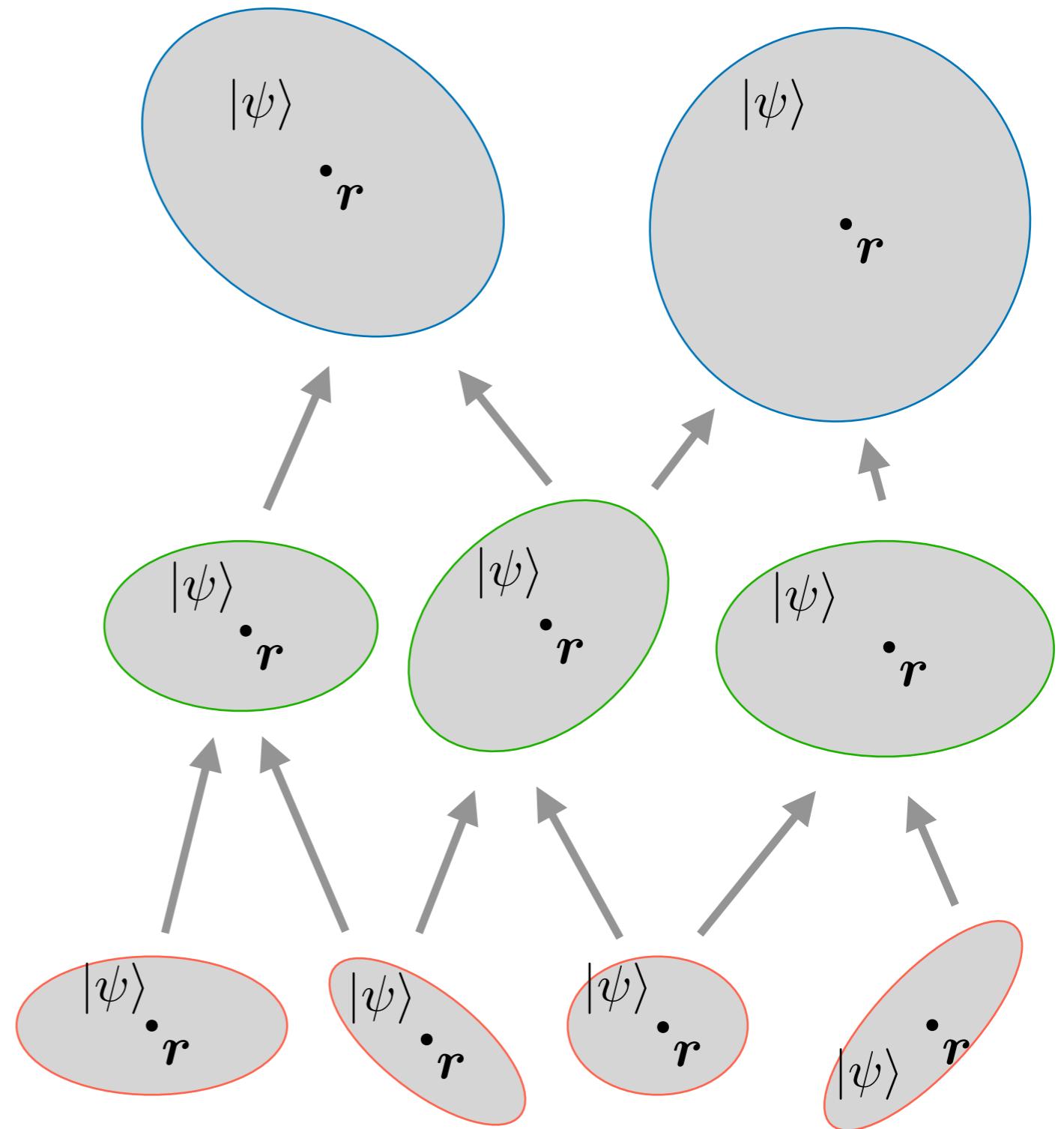
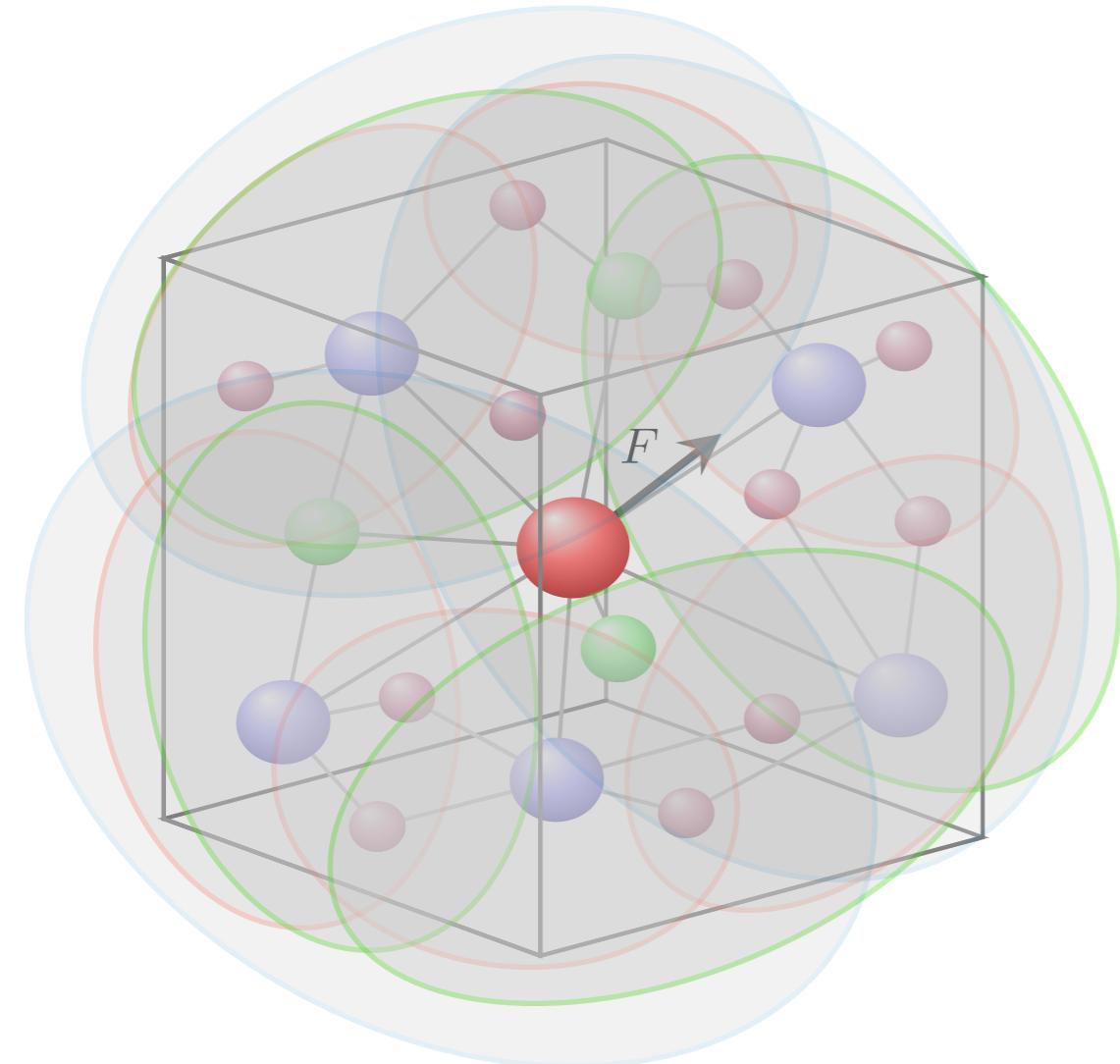
134k small organic molecules modeled with DFT  
[B3LYP/6-31G(2df,p)]

Targets: 15 physical/chemical properties

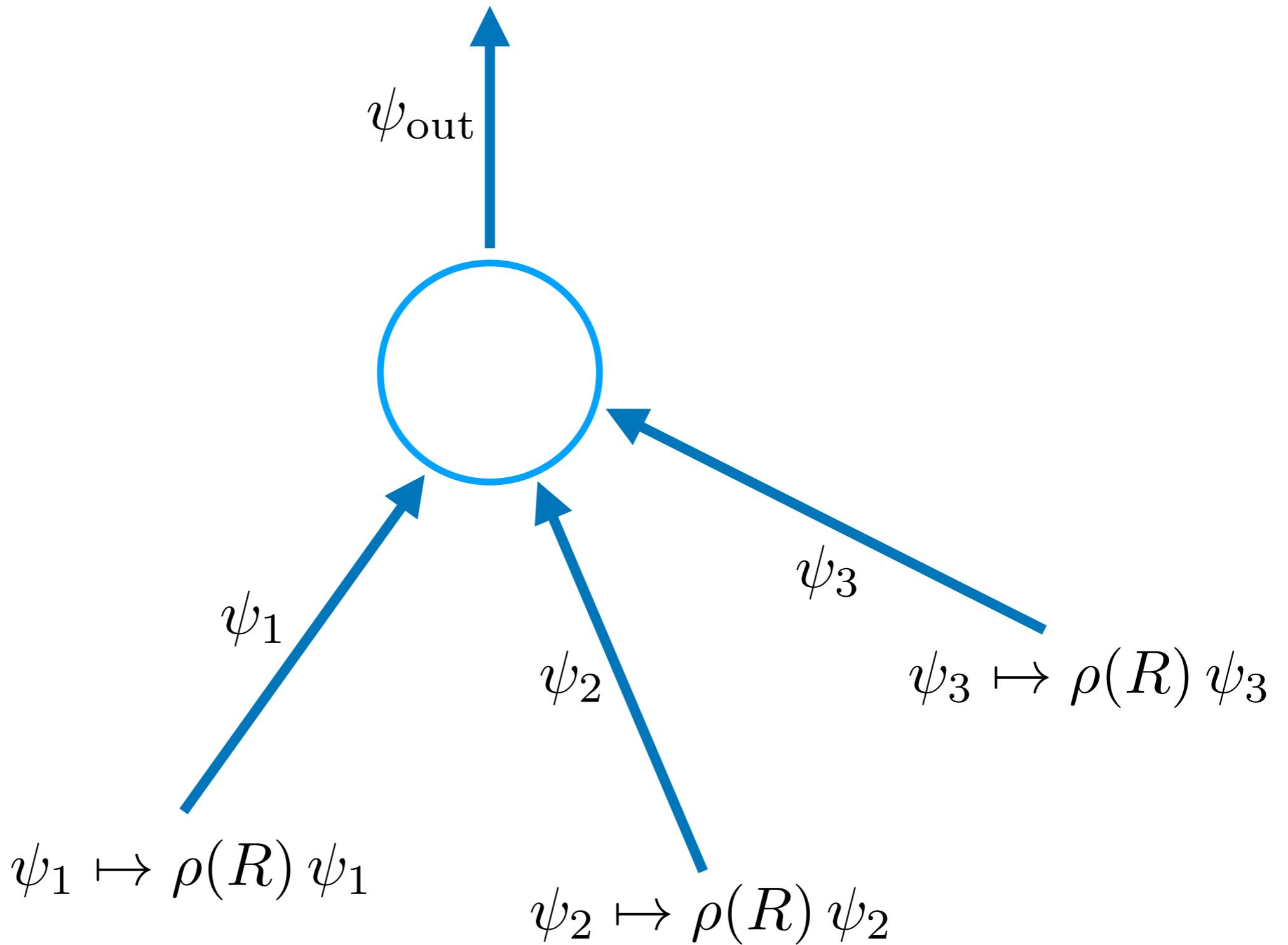


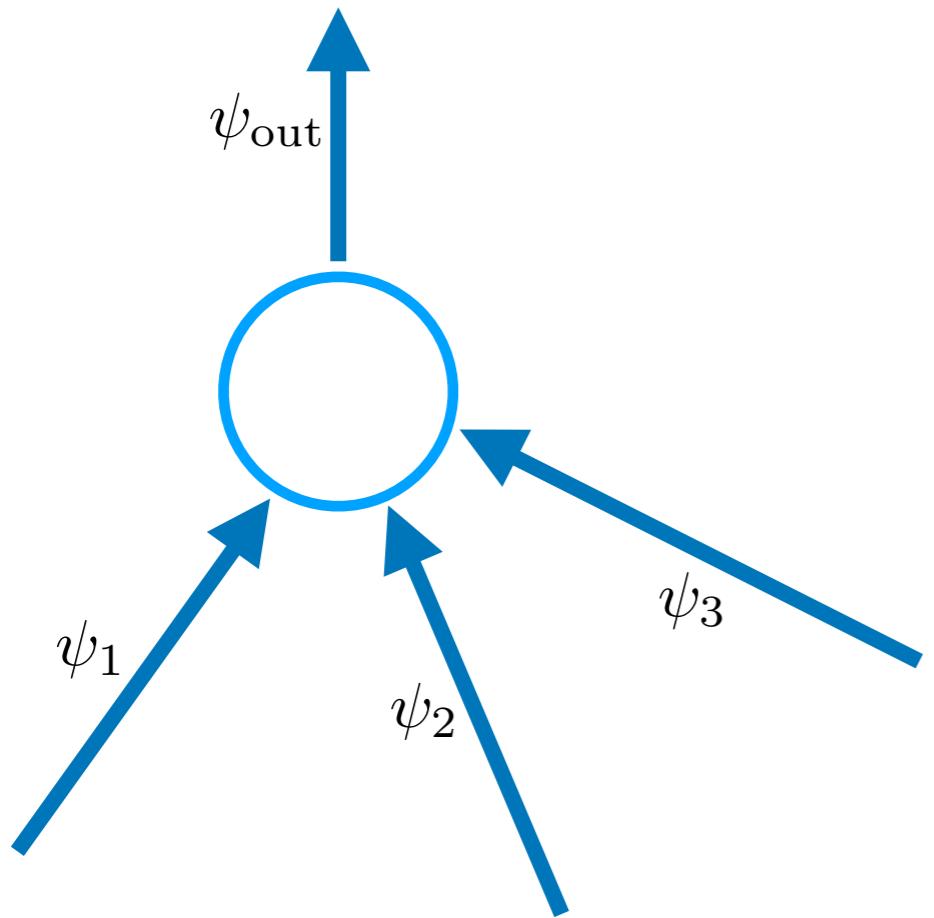
[Ramakrishnan et al., 2014] [von Lilienfeld et al.]

# Compositional structure



$$\psi_{\text{out}} \mapsto \rho(R) \psi_{\text{out}}$$





$$\psi_{\text{out}} = \sigma \left( \sum_i \psi_i * \chi \right)$$

$$\hat{\psi}_{\text{out}}^{\ell'} = \sigma \left( \sum_i \hat{\psi}_i^\ell \cdot \hat{\chi}^\ell \right)$$

But what form should the nonlinearity take?

# The Clebsch-Gordan product

The tensor product of two irreducible representations of a compact group  $G$  decomposes into irreducibles in the form

$$\rho_1(g) \otimes \rho_2(g) = \bigoplus_{\ell} \bigoplus_{m=1}^{\kappa(\ell)} C_m^\ell \cdot \rho_\ell(g) \cdot (C_m^\ell)^\dagger$$

# Clebsch-Gordan gates

General form of aggregation:

$$|\psi\rangle = \sigma\left(W \sum_i (\mathbf{r}_i - \mathbf{r})^{\otimes k} \otimes |\psi_i\rangle^{\otimes p}\right)$$

The nonlinearity is a low order Clebsch-Gordan product

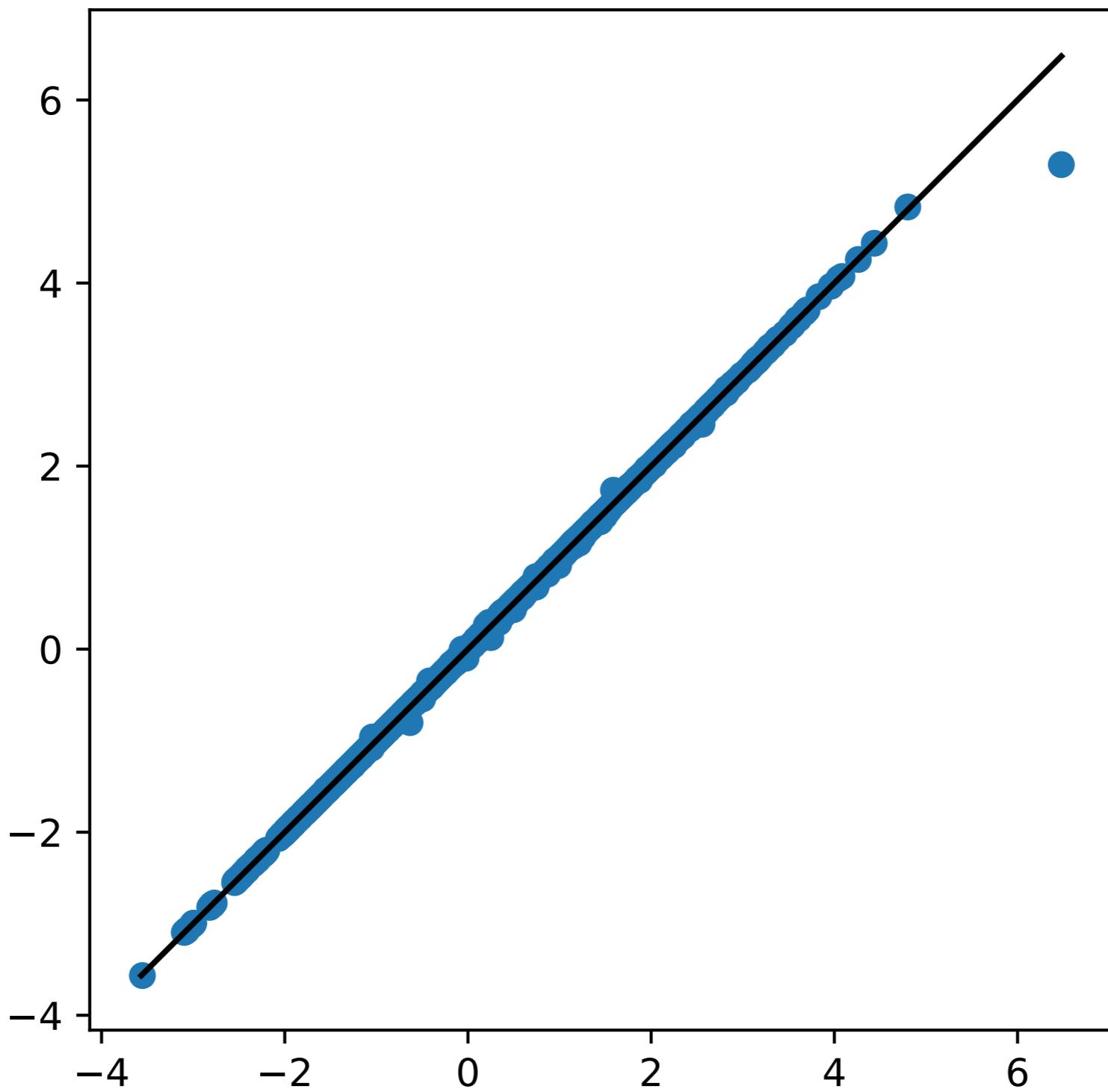
$$|\psi\rangle \mapsto \rho(g) |\psi\rangle$$

[“Tensor Field Networks” by Thomas et al., 2018] [“Clebsch-Gordan networks”, K., Trivedi & Zhen , 2018] [“3D steerable CNNs”, Weiler et al., 2018]

# QM9 results

	us	SchNett	google1	google2	googleML	L-Scat	T-Scat	NMP	CM
<b>alpha (bohr3)</b>	0.096	0.235	0.161	0.232	0.175	0.520	0.160	* 0.092	0.430
<b>Cv (cal/molK)</b>	* 0.031	0.033	0.084	0.097	0.044	0.100	0.049	0.040	0.120
<b>gap (eV)</b>	* 0.061	0.063	0.088	0.087	0.107	0.304	0.118	0.069	0.229
<b>homo (eV)</b>	* 0.036	0.041	0.057	0.055	0.066	0.177	0.085	0.043	0.133
<b>lumo (eV)</b>	* 0.032	0.034	0.063	0.062	0.084	0.234	0.076	0.038	0.183
<b>mu (D)</b>	0.046	0.033	0.247	0.101	0.334	0.630	0.340	* 0.030	0.450
<b>omega1 (cm-1)</b>	* 2.002	--	6.220	4.760	2.710	--	--	--	--
<b>r2 (bohr2)</b>	0.841	* 0.073	--	--	--	6.670	0.410	0.180	3.390
<b>U0 (eV)</b>	0.028	* 0.014	0.042	0.150	0.025	0.082	0.022	0.020	0.128
<b>zpve (meV)</b>	2.162	1.700	4.310	9.660	1.910	4.000	2.000	* 1.500	4.800

[Schütt et al., 2018] [Faber et al., 2017] [Eickenberg et al., 2018] [Zhang et al, 2018]



# MD-17 results (kcal/mol)

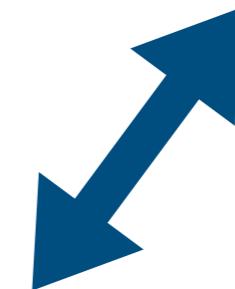
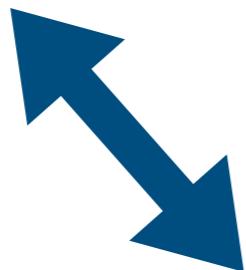
	NBody50k	Deep95k	DTNN50k	SchE50k	SchEF50k	GDML1k	SchE1k	SchEF1k
<b>Aspirin</b>	* 0.103	0.201	--	0.250	0.120	0.270	4.200	0.370
<b>Benzene</b>	* 0.035	0.065	0.040	0.080	0.070	0.070	1.190	0.080
<b>Ethanol</b>	* 0.029	0.055	--	0.070	0.050	0.150	0.930	0.080
<b>Malonaldehyde</b>	* 0.056	0.092	0.190	0.130	0.080	0.160	2.030	0.130
<b>Naphthalene</b>	* 0.043	0.095	--	0.200	0.110	0.120	3.580	0.160
<b>Salicylic_acid</b>	* 0.072	0.106	0.410	0.250	0.100	0.120	3.270	0.200
<b>Toluene</b>	* 0.042	0.085	0.180	0.160	0.090	0.120	2.950	0.120
<b>Uracil</b>	* 0.045	0.085	--	0.130	0.100	0.110	2.260	0.140



Theory

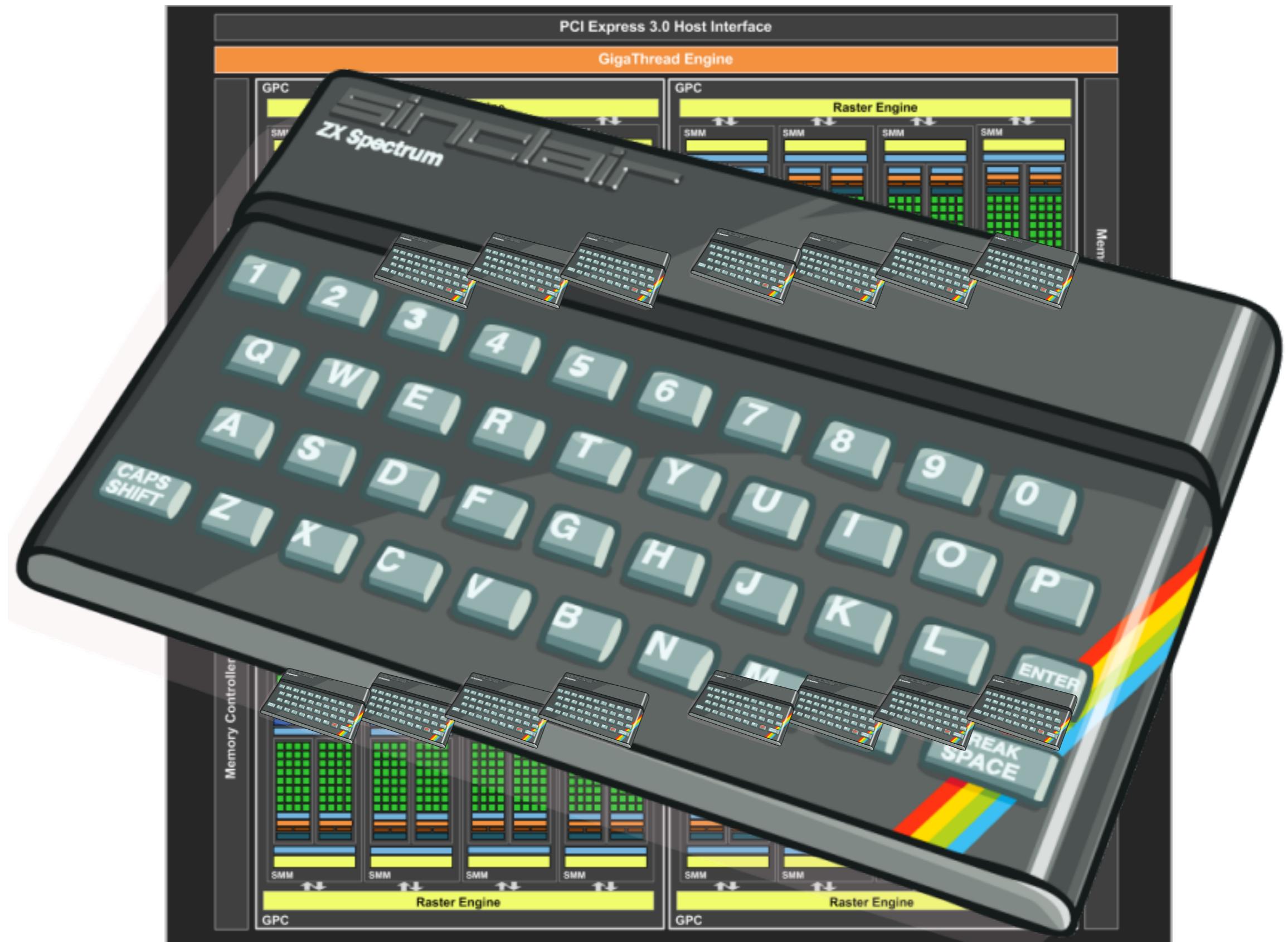


Experiment

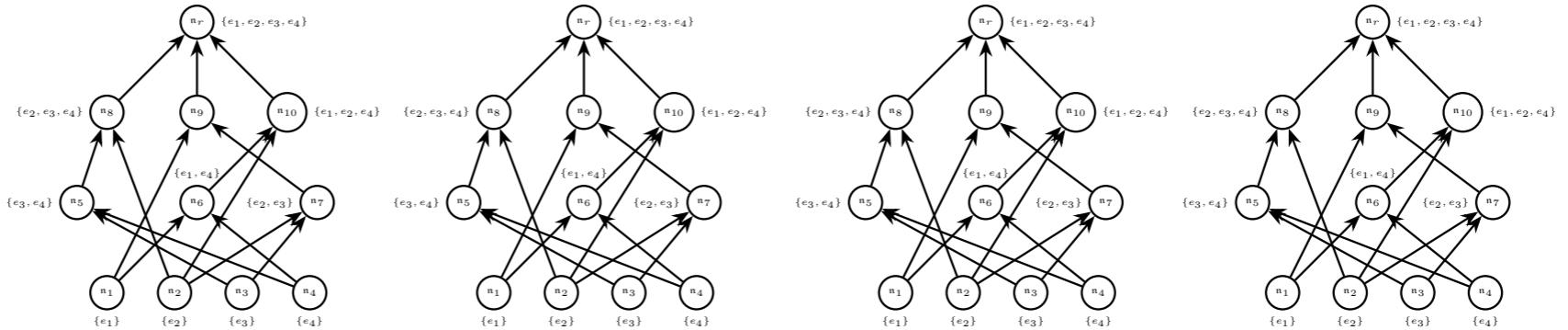


Implementation





# Computation graph



## Bytecode

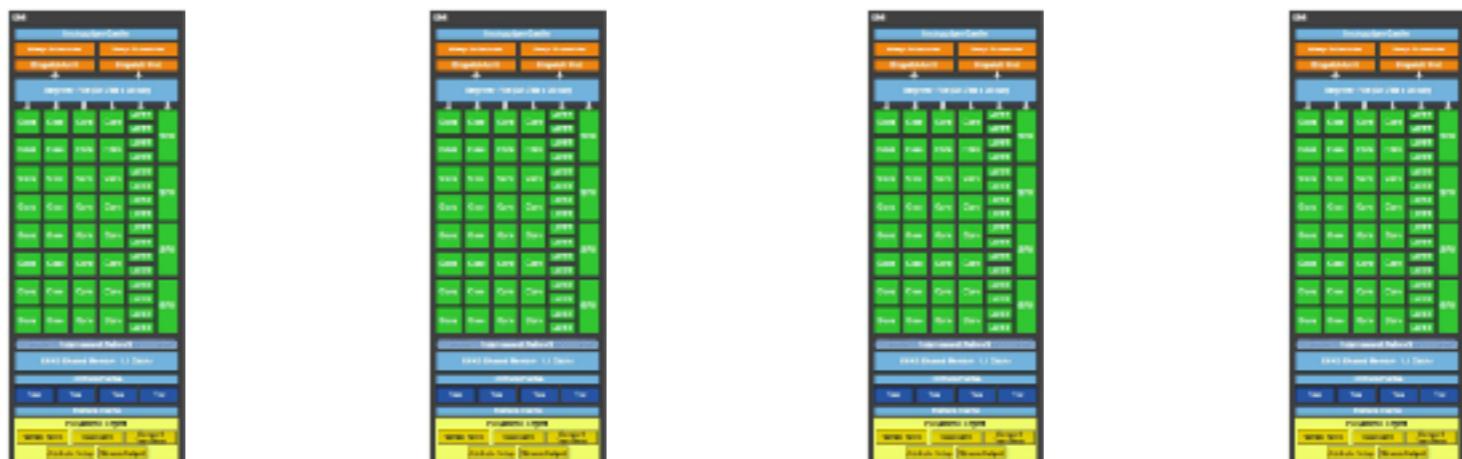
```
TPprogram CGproduct(){
    TPpart0 (1=0)[0m (n=1){
        input(0,0);
    }
    TPpart1 (1=1) (n=1){
        input(0,1);
    }
    TPpart2 (1=0) (n=1){
        input(1,0);
    }
    TPpart3 (1=1) (n=1){
        input(1,1);
    }
    TPpart4 (1=0) (n=1){
        input(2,0);
    }
    TPpart5 (1=1) (n=1){
        input(2,1);
    }
    TPpart6 (1=0) (n=2){
        CG(2,4)[0];
        CG(3,5)[1];
    }
    TPpart7 (1=1) (n=3){
        CG(2,5)[0];
        CG(3,4)[1];
        CG(3,5)[2];
    }
}
```

```
TPprogram CGproduct(){
    TPpart0 (1=0)[0m (n=1){
        input(0,0);
    }
    TPpart1 (1=1) (n=1){
        input(0,1);
    }
    TPpart2 (1=0) (n=1){
        input(1,0);
    }
    TPpart3 (1=1) (n=1){
        input(1,1);
    }
    TPpart4 (1=0) (n=1){
        input(2,0);
    }
    TPpart5 (1=1) (n=1){
        input(2,1);
    }
    TPpart6 (1=0) (n=2){
        CG(2,4)[0];
        CG(3,5)[1];
    }
    TPpart7 (1=1) (n=3){
        CG(2,5)[0];
        CG(3,4)[1];
        CG(3,5)[2];
    }
}
```

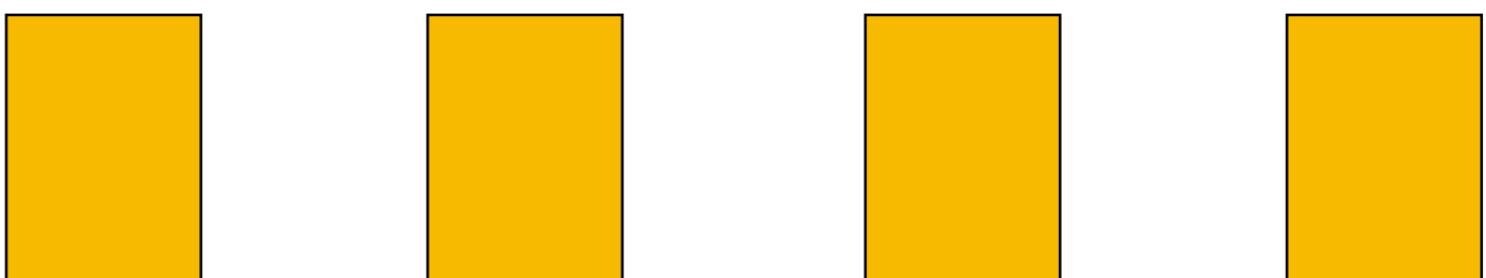
```
TPprogram CGproduct(){
    TPpart0 (1=0)[0m (n=1){
        input(0,0);
    }
    TPpart1 (1=1) (n=1){
        input(0,1);
    }
    TPpart2 (1=0) (n=1){
        input(1,0);
    }
    TPpart3 (1=1) (n=1){
        input(1,1);
    }
    TPpart4 (1=0) (n=1){
        input(2,0);
    }
    TPpart5 (1=1) (n=1){
        input(2,1);
    }
    TPpart6 (1=0) (n=2){
        CG(2,4)[0];
        CG(3,5)[1];
    }
    TPpart7 (1=1) (n=3){
        CG(2,5)[0];
        CG(3,4)[1];
        CG(3,5)[2];
    }
}
```

```
TPprogram CGproduct(){
    TPpart0 (1=0)[0m (n=1){
        input(0,0);
    }
    TPpart1 (1=1) (n=1){
        input(0,1);
    }
    TPpart2 (1=0) (n=1){
        input(1,0);
    }
    TPpart3 (1=1) (n=1){
        input(1,1);
    }
    TPpart4 (1=0) (n=1){
        input(2,0);
    }
    TPpart5 (1=1) (n=1){
        input(2,1);
    }
    TPpart6 (1=0) (n=2){
        CG(2,4)[0];
        CG(3,5)[1];
    }
    TPpart7 (1=1) (n=3){
        CG(2,5)[0];
        CG(3,4)[1];
        CG(3,5)[2];
    }
}
```

Same bytecode  
interpreter on each  
streaming  
multiprocessor



Result



```

VertexFn inputs;

VertexFn L1 =
aggregate1D(perc(inputs,W1));
EdgeFn L2 =
scatter(L1)*scatter(L1);
EdgeFn L3 =
perc<ReLU>(scatter(L1),W2);
EdgeFn L4 =
perc<ReLU>(concat(L2,L3));
VertexFn L5 = gather(L4);
VertexFn L6 =
perc<sigmoid>(L5,W4);

GlobalFn L8 =
average(aggregate0D(L6));
GlobalFn L9 =
perc<sigmoid>(L8,W5);
GlobalFn
output=squared_error_loss(L9,
target);

```

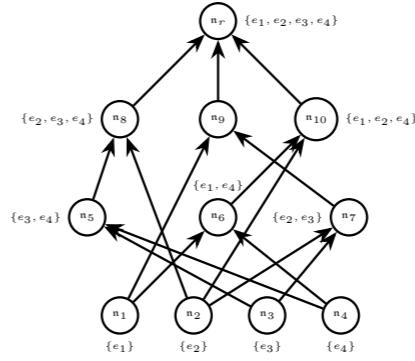
## GNN grammar

```

CCNlayer0D 3
CCNlayer1D 3
ReLUlayer 2
CCNlayer0D 4
ReLUlayer 3
AveragingLayer 3
LinearFClayer 1
SquaredLossLayer
1

```

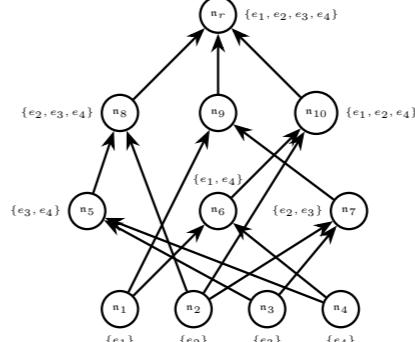
## GNNlayerDescr



```

TPprogram CGproduct(){
TPpart0 (1=0)[0m (n=1){
    input(0,0);
}
TPpart1 (1=1) (n=1){
    input(0,1);
}
TPpart2 (1=0) (n=1){
    input(1,0);
}
TPpart3 (1=1) (n=1){
    input(1,1);
}
TPpart4 (1=0) (n=1){
    input(2,0);
}
TPpart5 (1=1) (n=1){
    input(2,1);
}
TPpart6 (1=0) (n=2){
    CG(2,4)[0];
    CG(3,5)[1];
}
TPpart7 (1=1) (n=3){
    CG(2,5)[0];
    CG(3,4)[1];
    CG(3,5)[2];
}

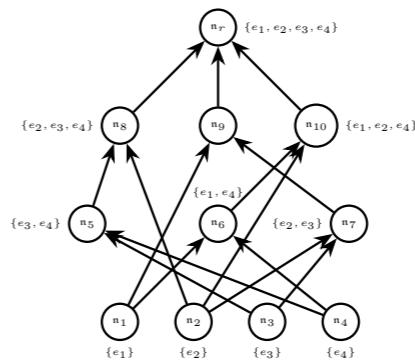
```



```

TPprogram CGproduct(){
TPpart0 (1=0)[0m (n=1){
    input(0,0);
}
TPpart1 (1=1) (n=1){
    input(0,1);
}
TPpart2 (1=0) (n=1){
    input(1,0);
}
TPpart3 (1=1) (n=1){
    input(1,1);
}
TPpart4 (1=0) (n=1){
    input(2,0);
}
TPpart5 (1=1) (n=1){
    input(2,1);
}
TPpart6 (1=0) (n=2){
    CG(2,4)[0];
    CG(3,5)[1];
}
TPpart7 (1=1) (n=3){
    CG(2,5)[0];
    CG(3,4)[1];
    CG(3,5)[2];
}

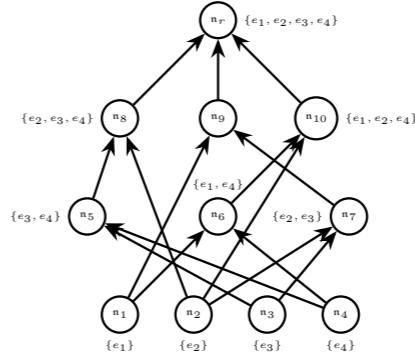
```



```

TPprogram CGproduct(){
TPpart0 (1=0)[0m (n=1){
    input(0,0);
}
TPpart1 (1=1) (n=1){
    input(0,1);
}
TPpart2 (1=0) (n=1){
    input(1,0);
}
TPpart3 (1=1) (n=1){
    input(1,1);
}
TPpart4 (1=0) (n=1){
    input(2,0);
}
TPpart5 (1=1) (n=1){
    input(2,1);
}
TPpart6 (1=0) (n=2){
    CG(2,4)[0];
    CG(3,5)[1];
}
TPpart7 (1=1) (n=3){
    CG(2,5)[0];
    CG(3,4)[1];
    CG(3,5)[2];
}

```

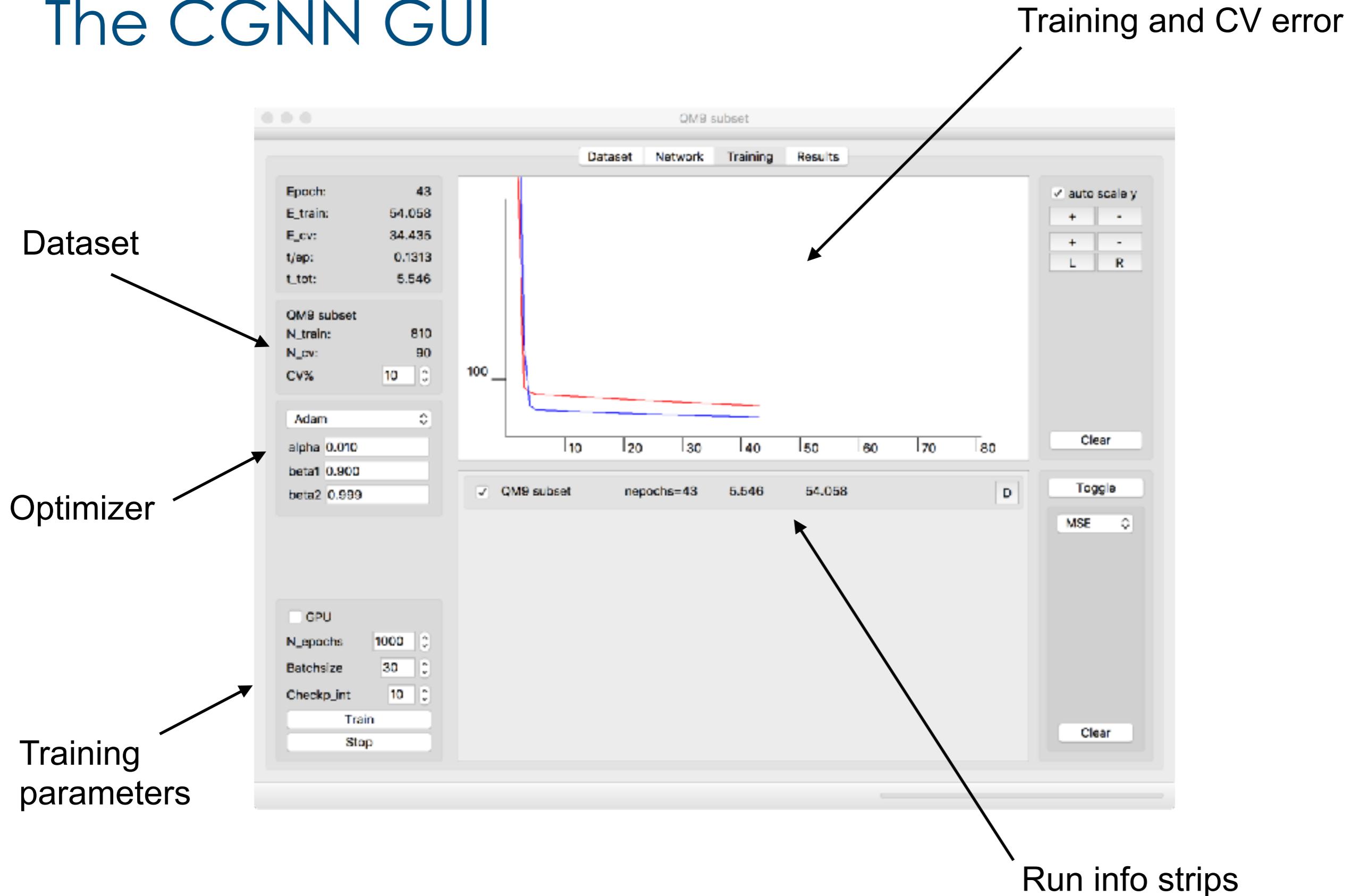


```

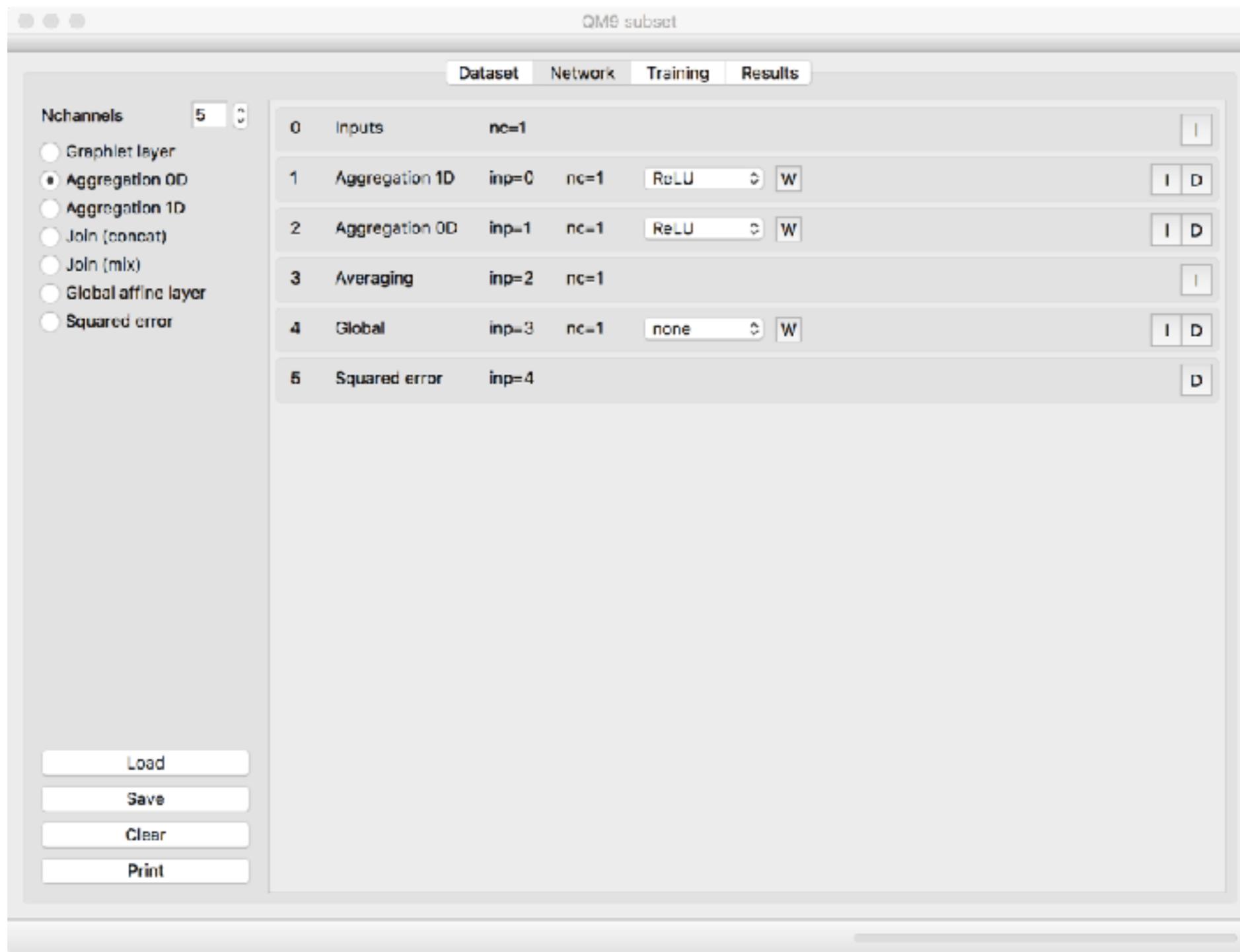
TPprogram CGproduct(){
TPpart0 (1=0)[0m (n=1){
    input(0,0);
}
TPpart1 (1=1) (n=1){
    input(0,1);
}
TPpart2 (1=0) (n=1){
    input(1,0);
}
TPpart3 (1=1) (n=1){
    input(1,1);
}
TPpart4 (1=0) (n=1){
    input(2,0);
}
TPpart5 (1=1) (n=1){
    input(2,1);
}
TPpart6 (1=0) (n=2){
    CG(2,4)[0];
    CG(3,5)[1];
}
TPpart7 (1=1) (n=3){
    CG(2,5)[0];
    CG(3,4)[1];
    CG(3,5)[2];
}

```

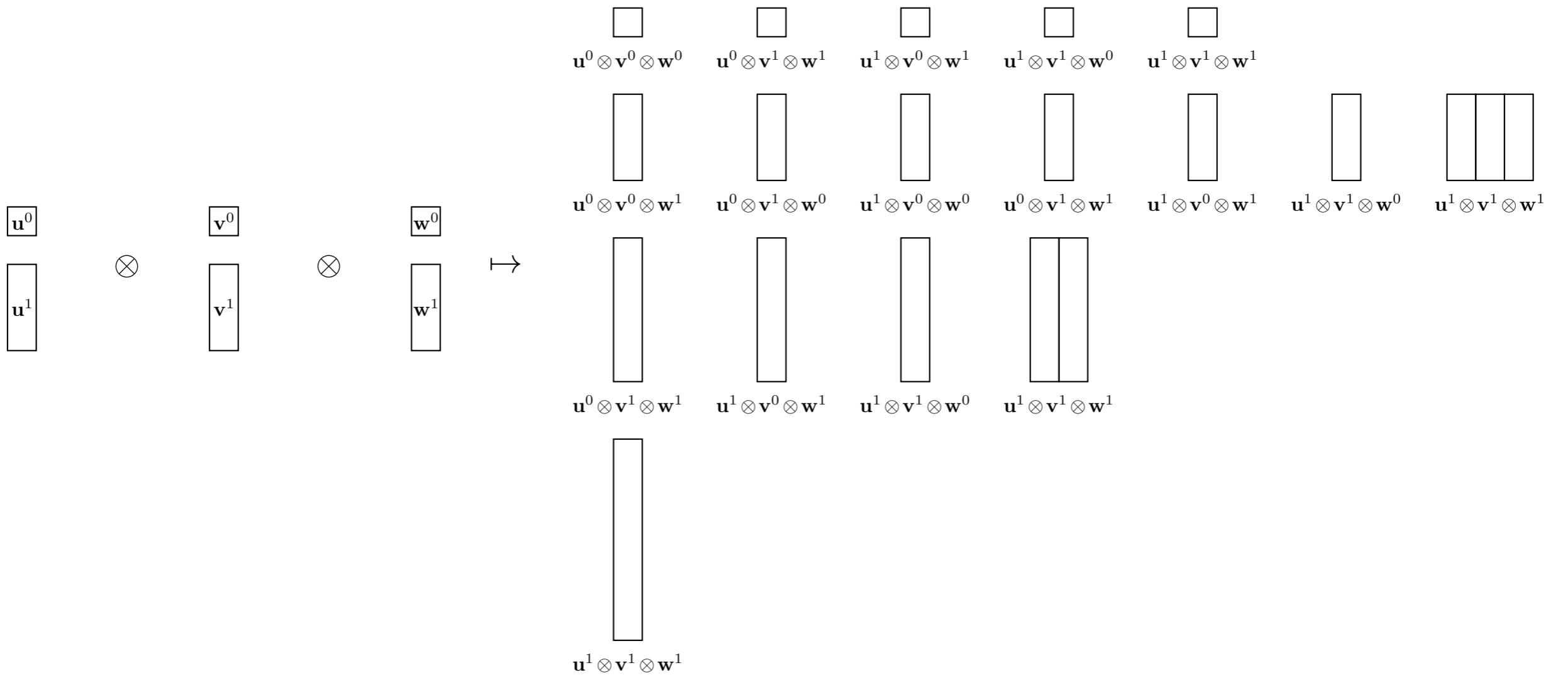
# The CGNN GUI

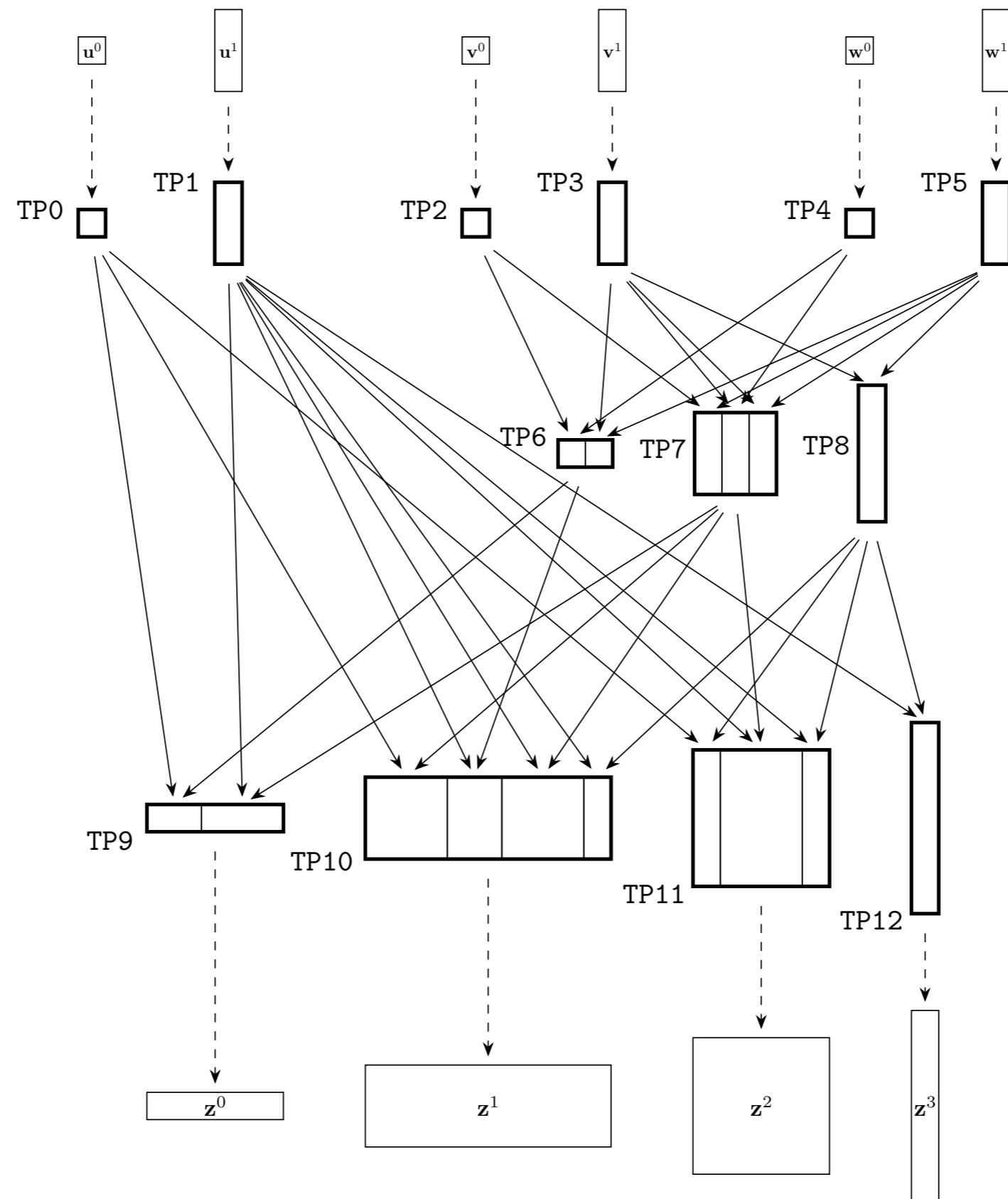


# The CGNN GUI



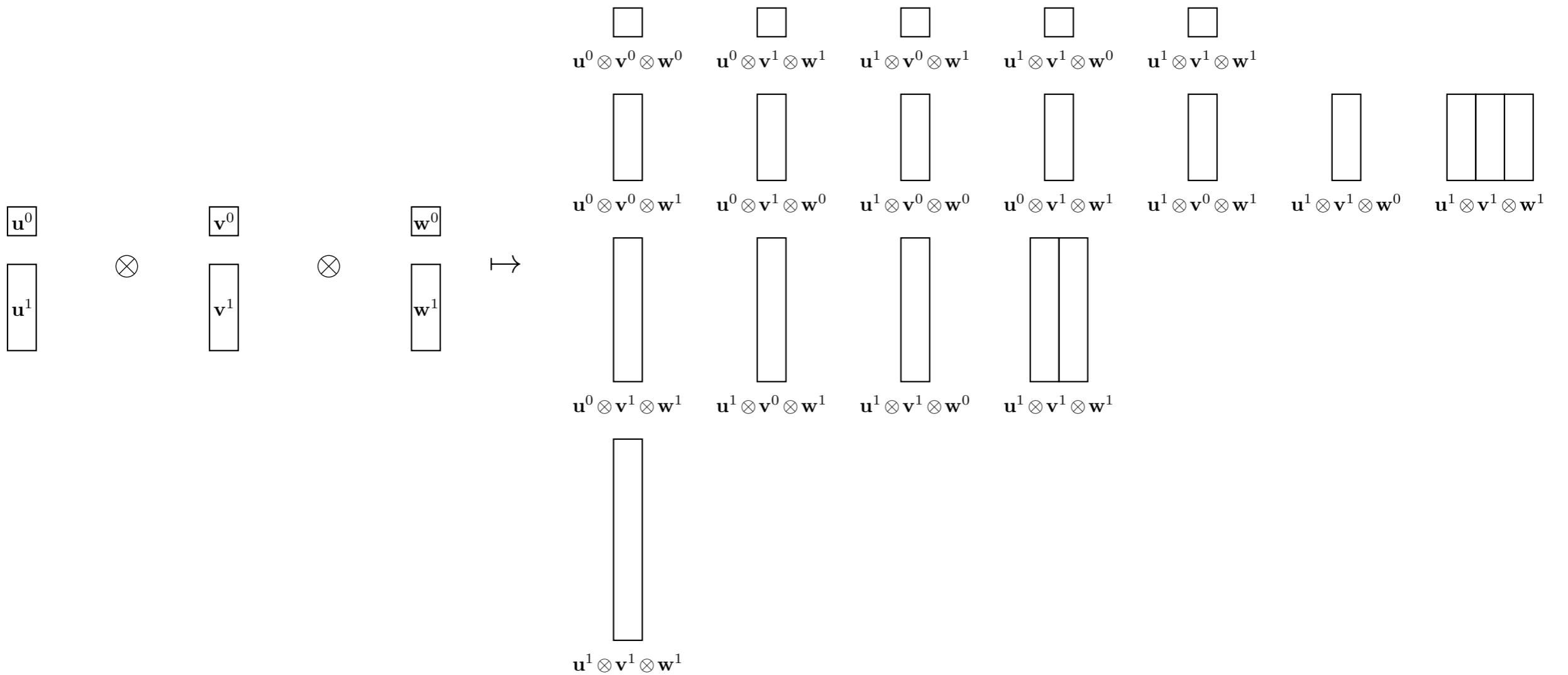
$$\begin{array}{c}
 \boxed{\mathbf{u}^1} \\
 \otimes \\
 \boxed{\mathbf{v}^1}
 \end{array}
 \rightarrow
 \begin{array}{c}
 \boxed{\mathbf{z}^0} \\
 \boxed{\mathbf{z}^1} \\
 \boxed{\mathbf{z}^2}
 \end{array}$$

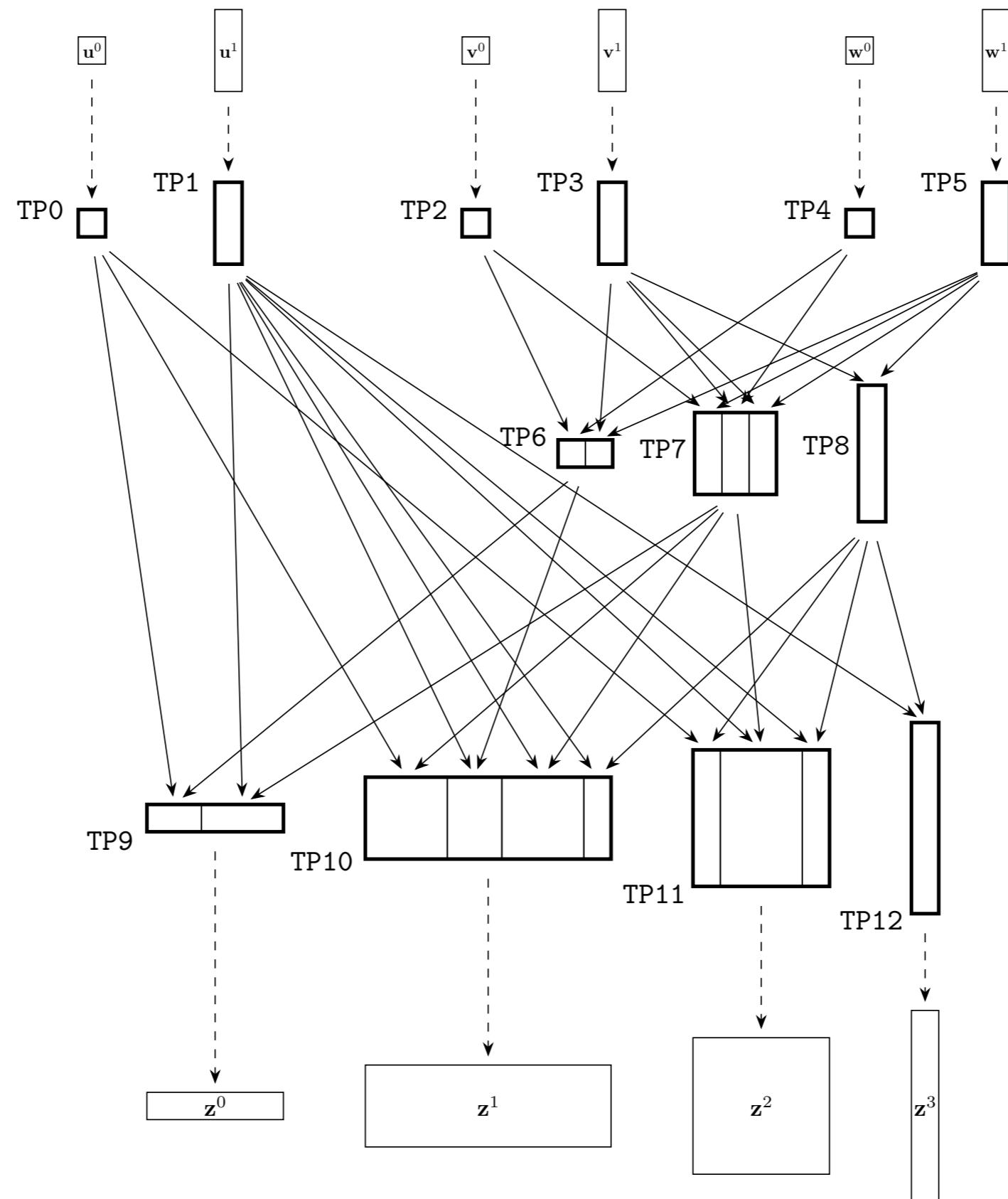




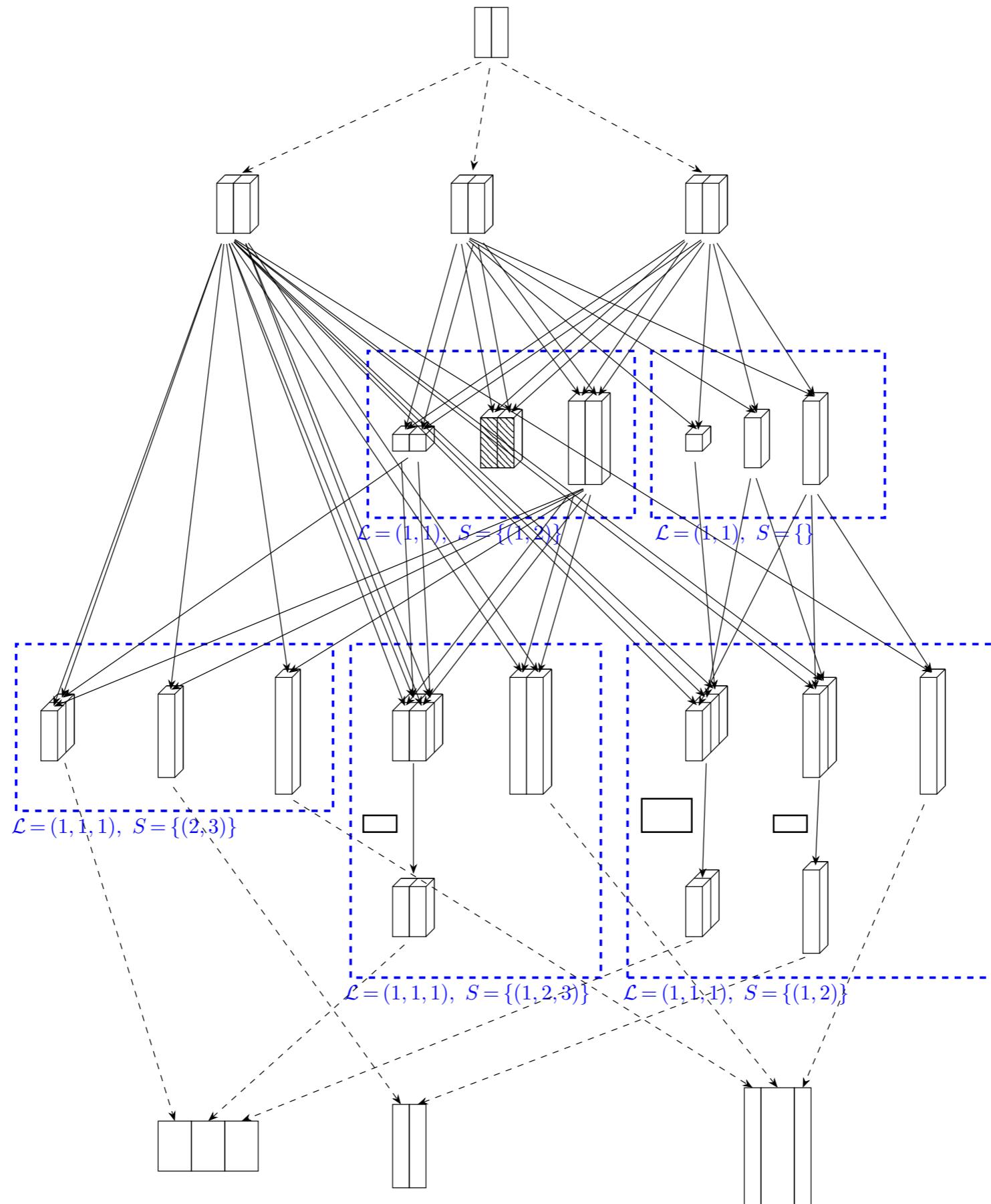
```
1  TPprogram  CGproduct(){
2      TPpart0 (l=0) [0m (n=1){
3          input(0,0);
4      }
5      TPpart1 (l=1) (n=1){
6          input(0,1);
7      }
8      TPpart2 (l=0) (n=1){
9          input(1,0);
10     }
11     TPpart3 (l=1) (n=1){
12         input(1,1);
13     }
14     TPpart4 (l=0) (n=1){
15         input(2,0);
16     }
17     TPpart5 (l=1) (n=1){
18         input(2,1);
19     }
20     TPpart6 (l=0) (n=2){
21         CG(2,4)[0];
22         CG(3,5)[1];
23     }
24     TPpart7 (l=1) (n=3){
25         CG(2,5)[0];
26         CG(3,4)[1];
27         CG(3,5)[2];
28     }
29     TPpart8 (l=2) (n=1){
30         CG(3,5)[0];
31     }
32     TPpart9 (l=0) (n=5){
33         output(0);
34         CG(0,6)[0];
35         CG(1,7)[2];
36     }
37     TPpart10 (l=1) (n=9){
38         output(1);
39         CG(0,7)[0];
40         CG(1,6)[3];
41         CG(1,7)[5];
42 }
```

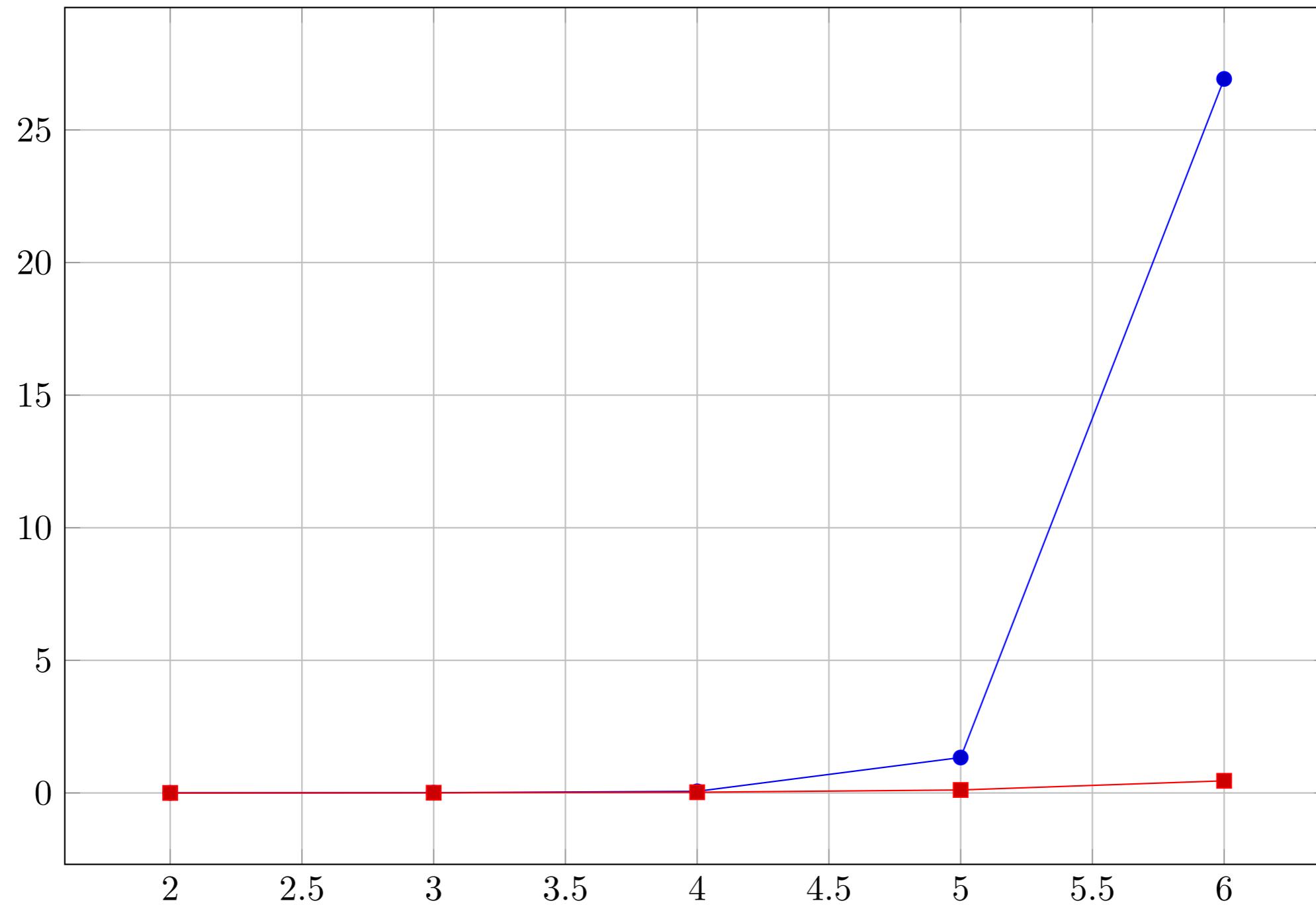
$$\begin{array}{c}
 \boxed{\mathbf{u}^1} \\
 \otimes \\
 \boxed{\mathbf{v}^1}
 \end{array}
 \rightarrow
 \begin{array}{c}
 \boxed{\mathbf{z}^0} \\
 \boxed{\mathbf{z}^1} \\
 \boxed{\mathbf{z}^2}
 \end{array}$$





$$\begin{array}{c} \boxed{\mathbf{v}^1} \\ \otimes \\ \boxed{\mathbf{v}^1} \\ \mapsto \\ \begin{array}{c} \boxed{\mathbf{w}^0} \\ \boxed{\mathbf{w}^1} \\ \boxed{\mathbf{w}^2} \end{array} \end{array}$$





## 1. Compositional structure

## 2. Covariance

→ Fourier space activations



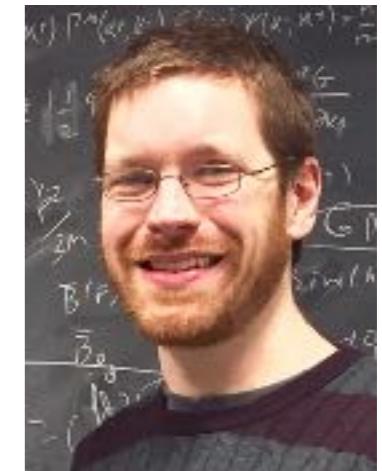
Shubhendu  
Trivedi



Hy Truong  
Son



Horace Pan



Brandon  
Anderson