

Multiplicative interactions

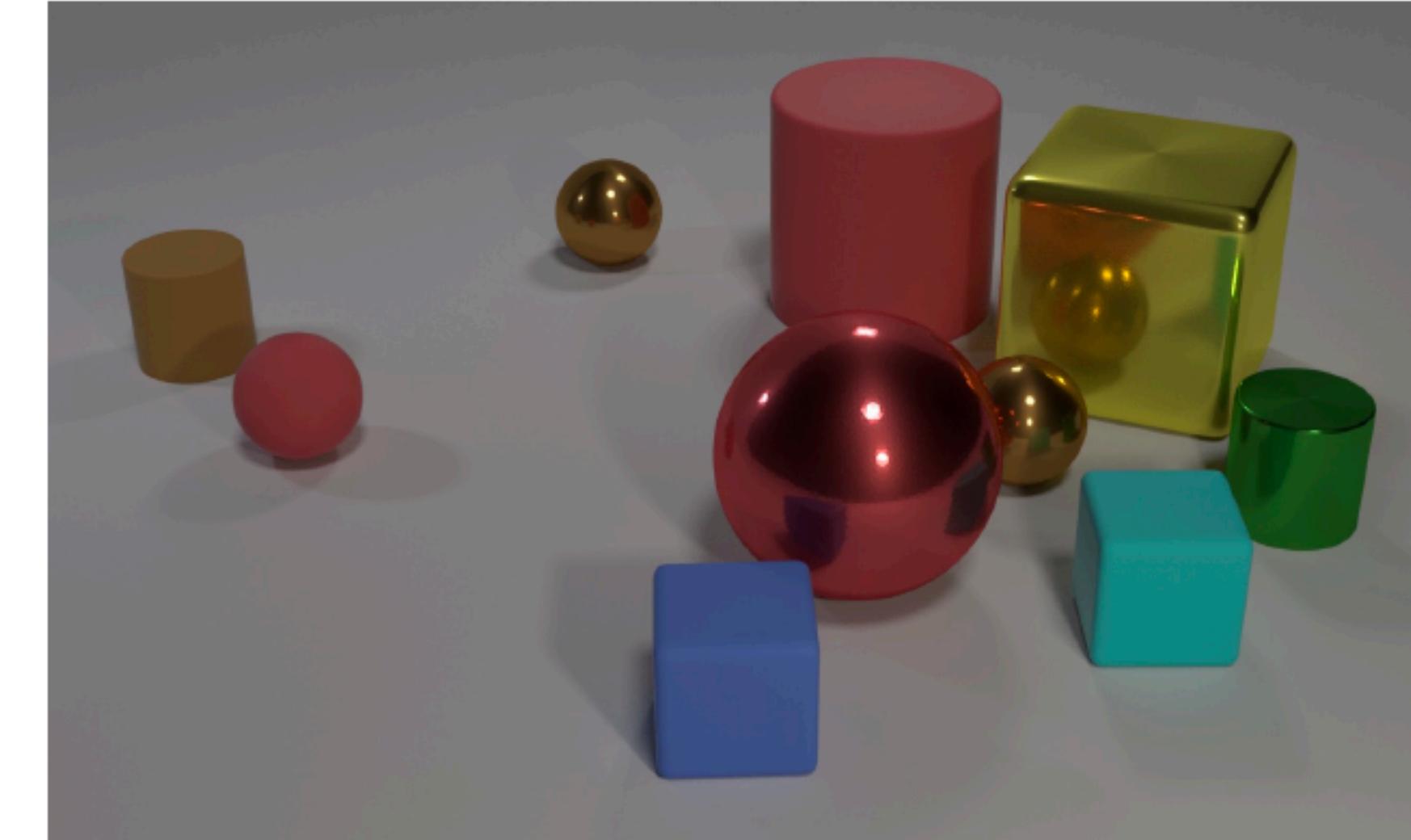
Alex Fedorov

January 22, 2021

Motivation

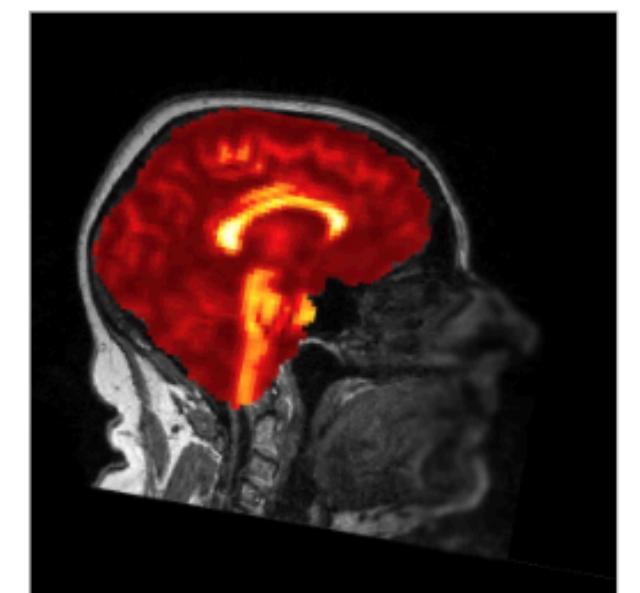
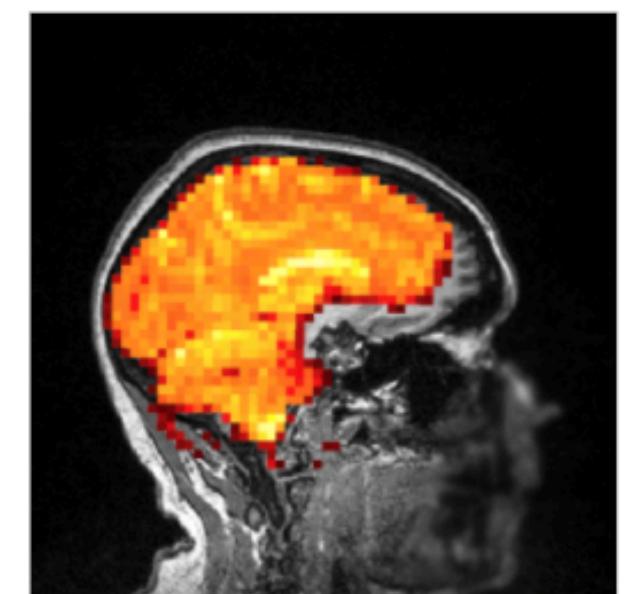
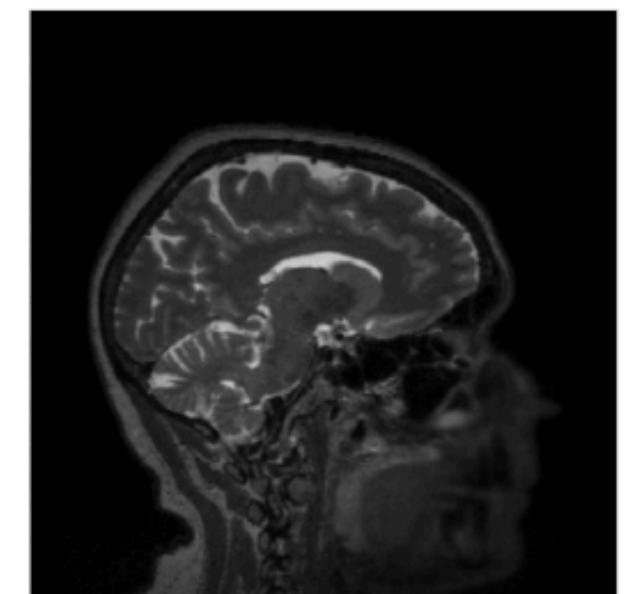
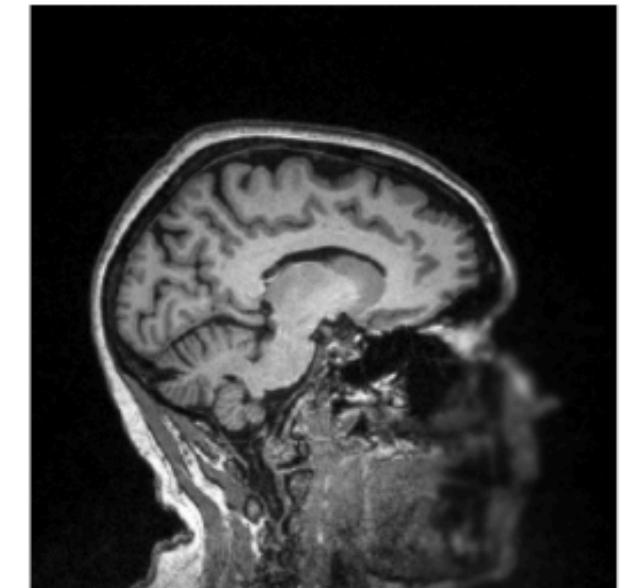


<https://weheartit.com/entry/280850620>

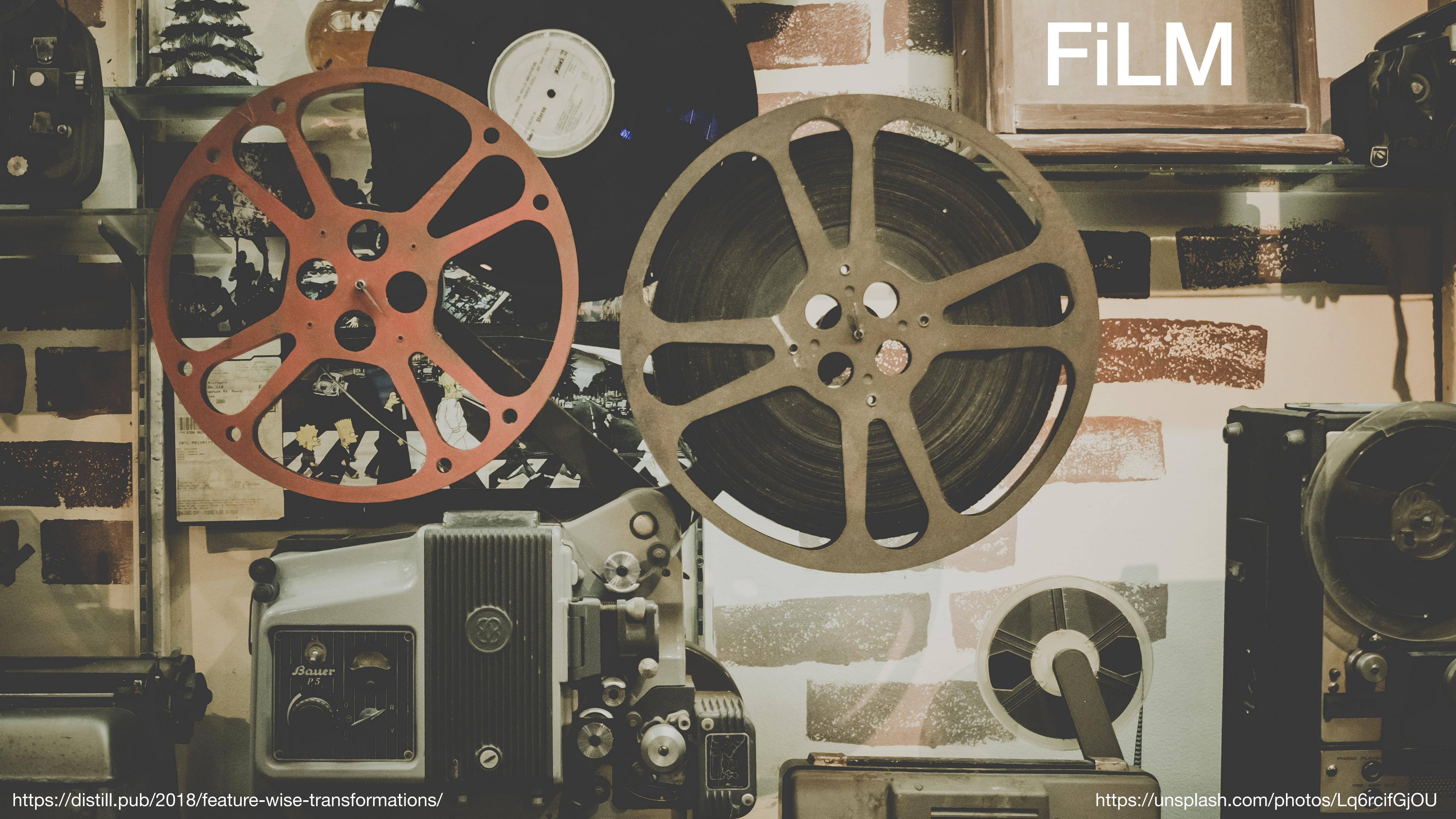


- Q:** Are there an equal number of large things and metal spheres?
Q: What size is the cylinder that is left of the brown metal thing that is left of the big sphere? **Q:** There is a sphere with the same size as the metal cube; is it made of the same material as the small red sphere?
Q: How many objects are either small cylinders or metal things?

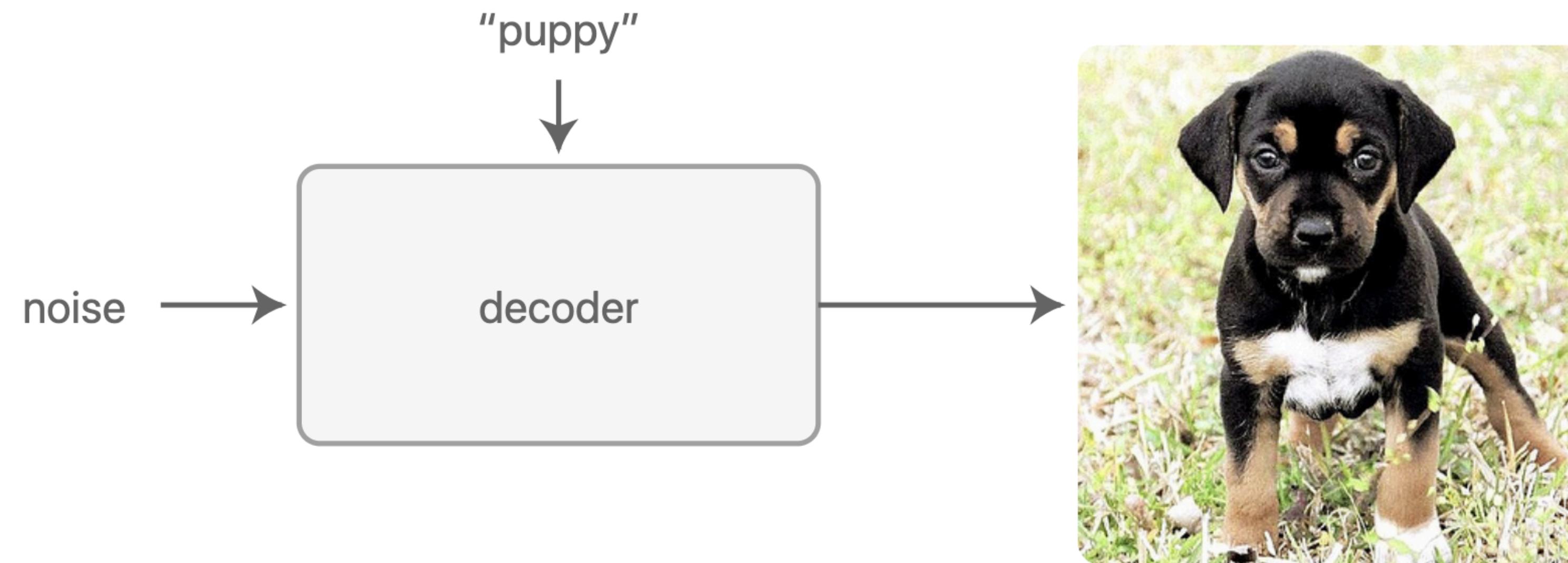
<https://towardsdatascience.com/deep-learning-and-visual-question-answering-c8c8093941bc>



FiLM



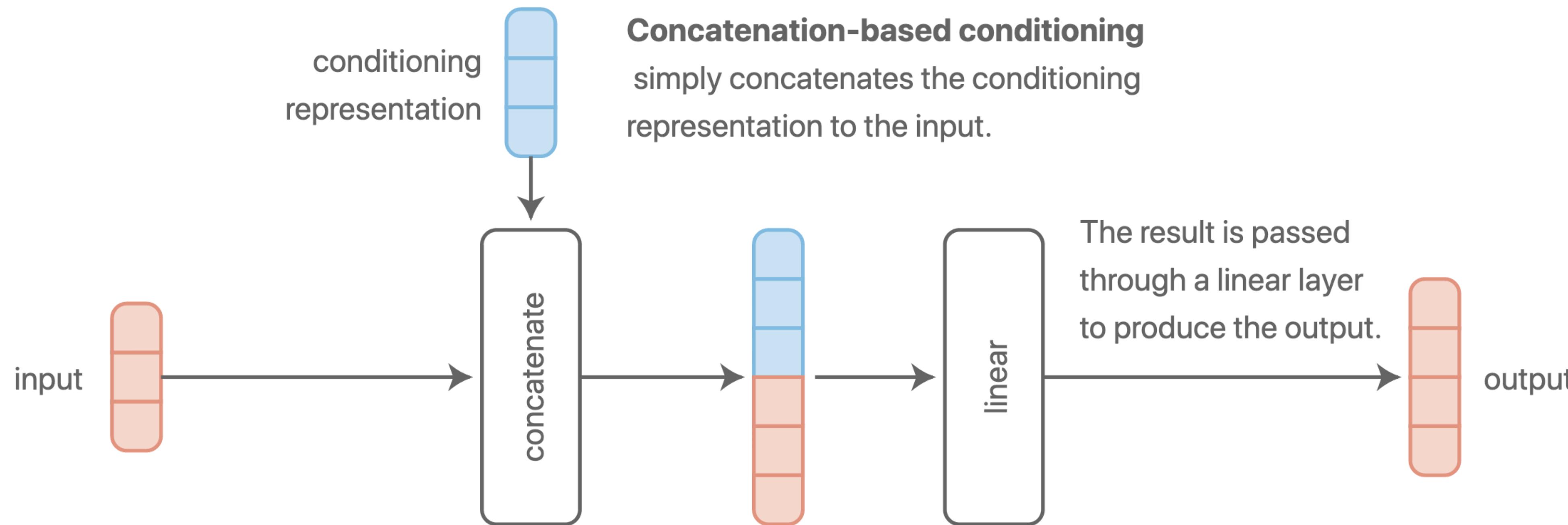
Conditioning



A **decoder-based generative model** maps a source of noise to a sample in the context of the "puppy" class.

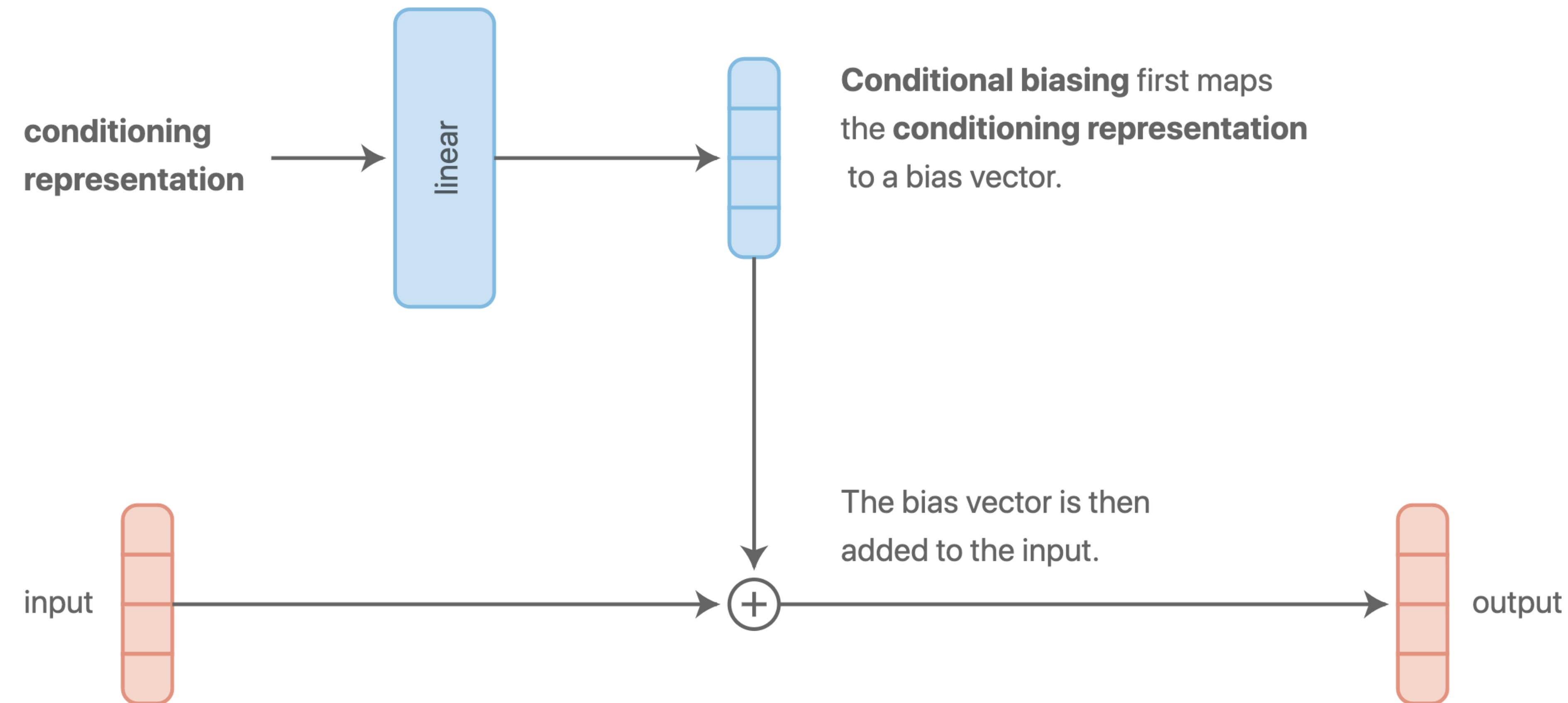
- What if we have 1000 classes?

Concatenation-based conditioning



- implicit approach
- feature aggregation / detection

Conditional biasing

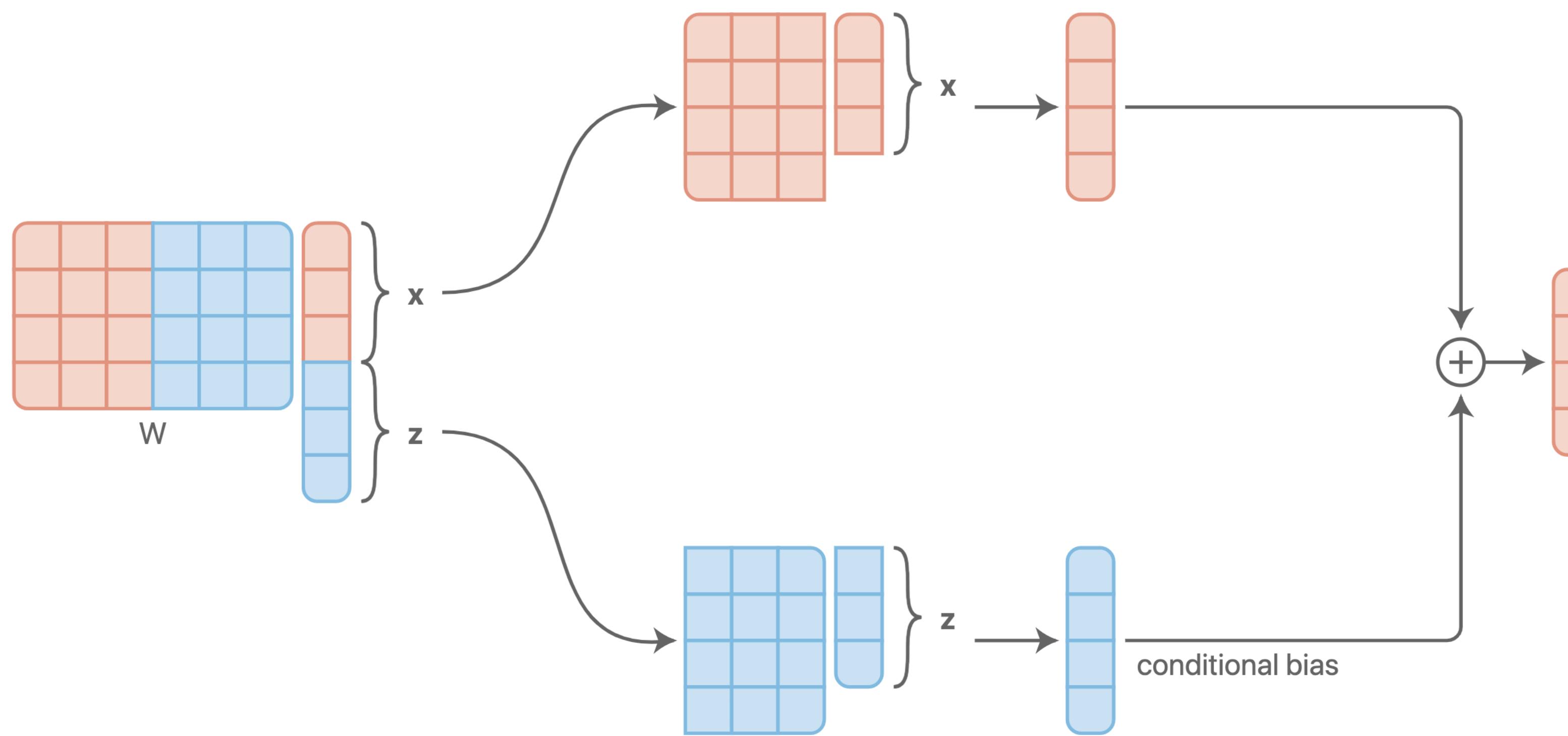


Concatenation == Biasing

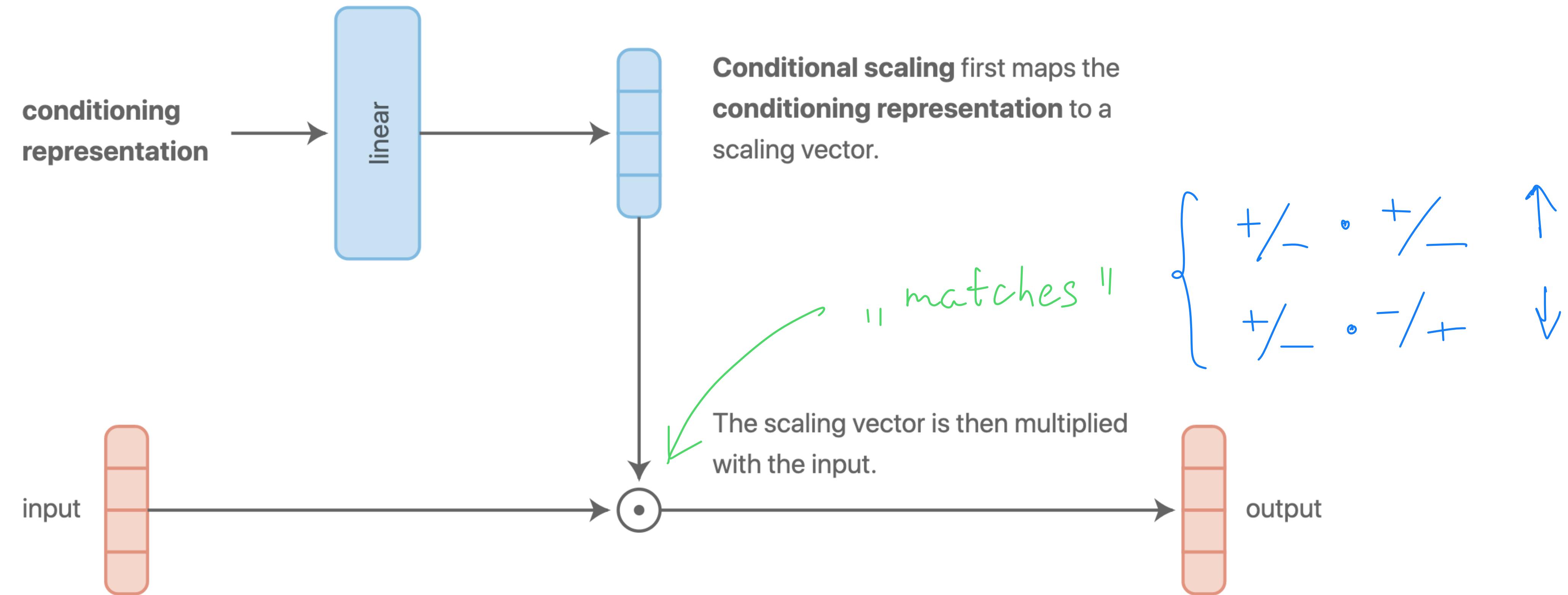
Concatenation-based conditioning
is equivalent to **conditional biasing**.

We can decompose the matrix-
vector product into two matrix-
vector subproducts.

We can then add the
resulting two vectors.
The \mathbf{z} -dependent vector
is a conditional bias.



Conditional scaling

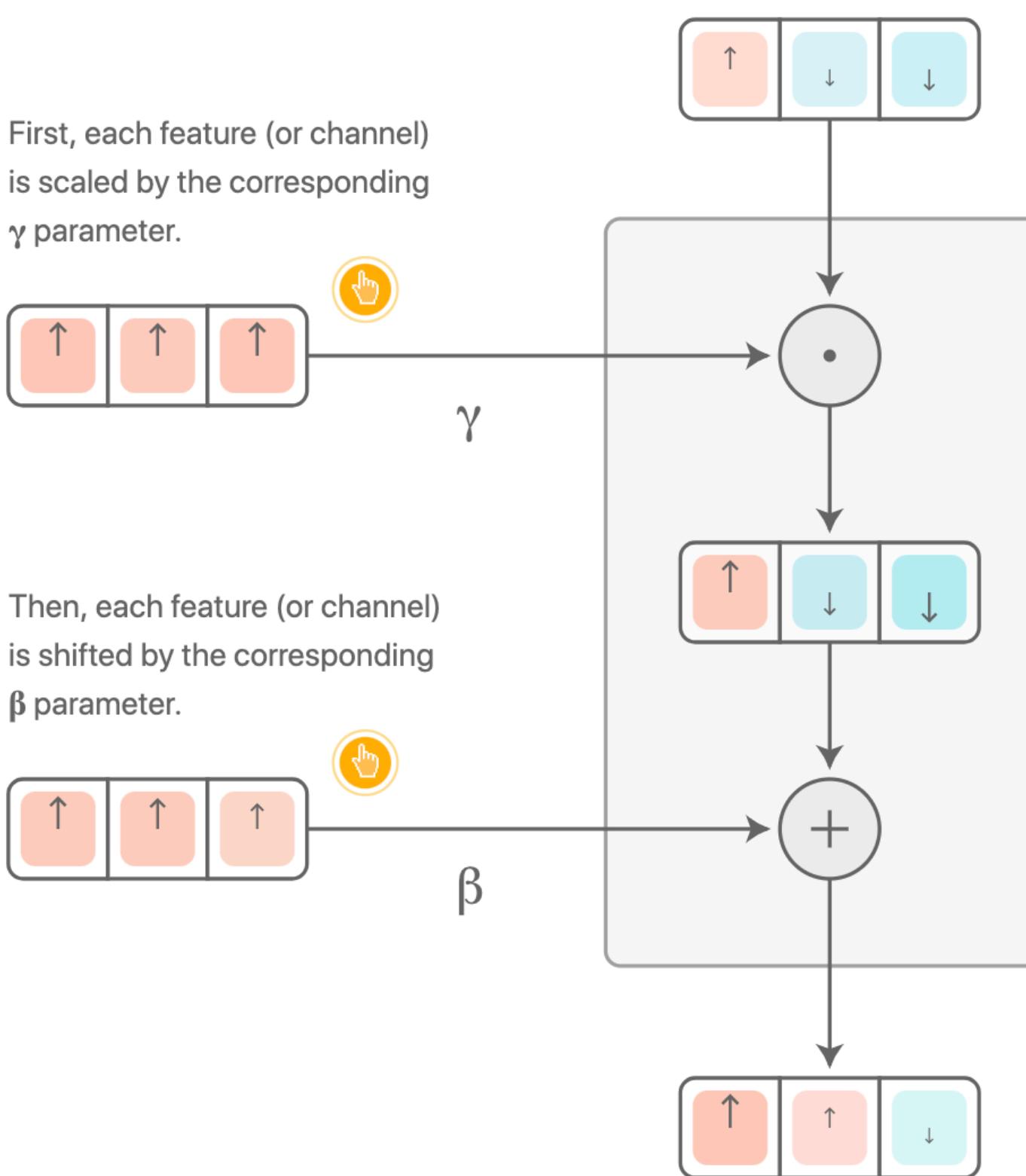


Example : feature wise sigmoidal gating

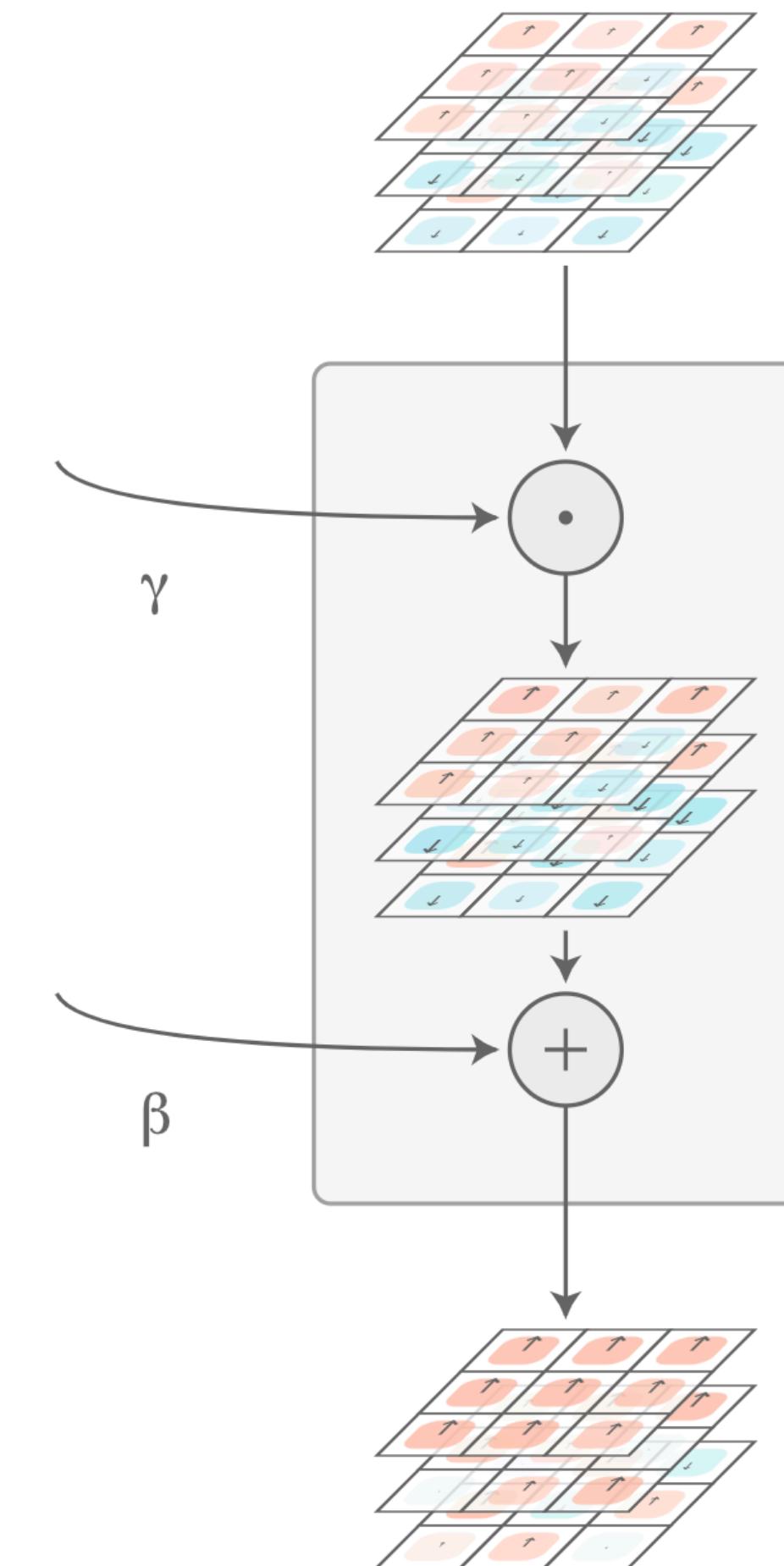
$$\text{FiLM}(\mathbf{x}) = \gamma(\mathbf{z}) \odot \mathbf{x} + \beta(\mathbf{z})$$

- Scales linearly $\mathcal{O}(N)$
number
of features
- Enough capacity for most scenarios

In a **fully-connected** network,
FiLM applies a different affine transformation to each feature.



In a **convolutional** network,
FiLM applies a different affine transformation to each channel, consistent across spatial locations.

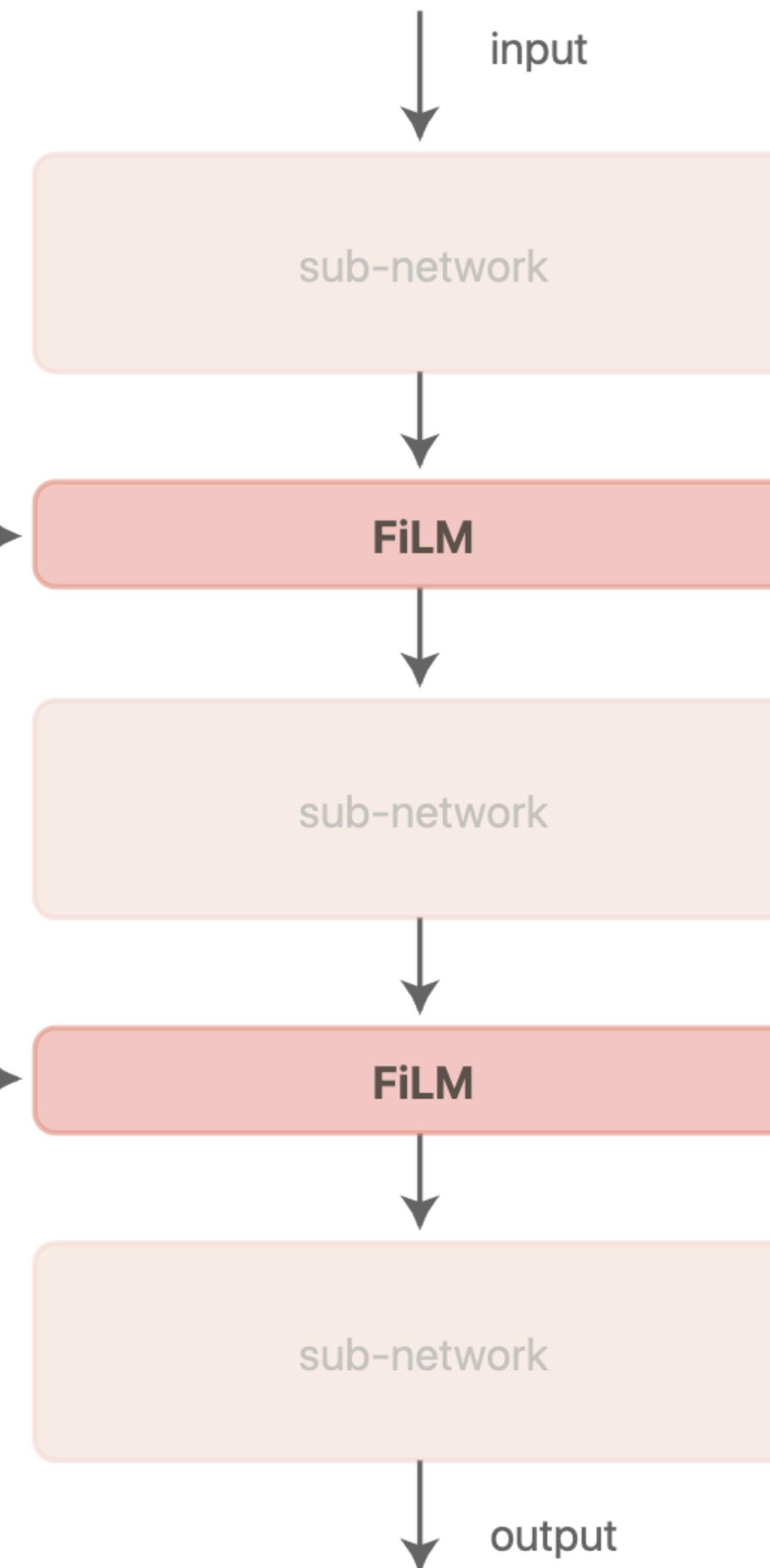


The **FiLM generator** processes the conditioning information and produces parameters that describe how the target network should alter its computation.



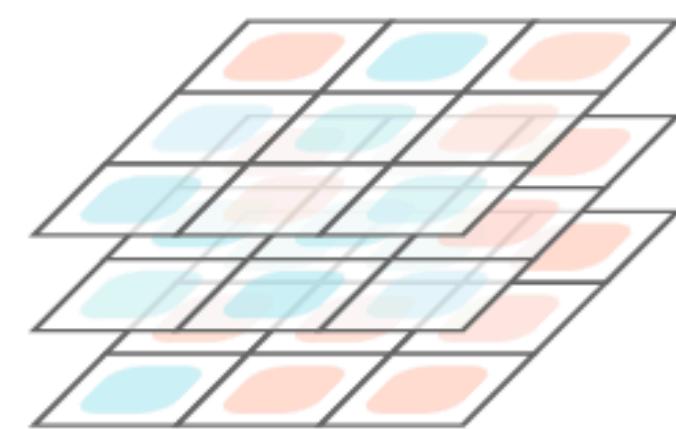
- Compounded across multiple layers

Here, the **FiLM-ed network's** computation is conditioned by two FiLM layers.

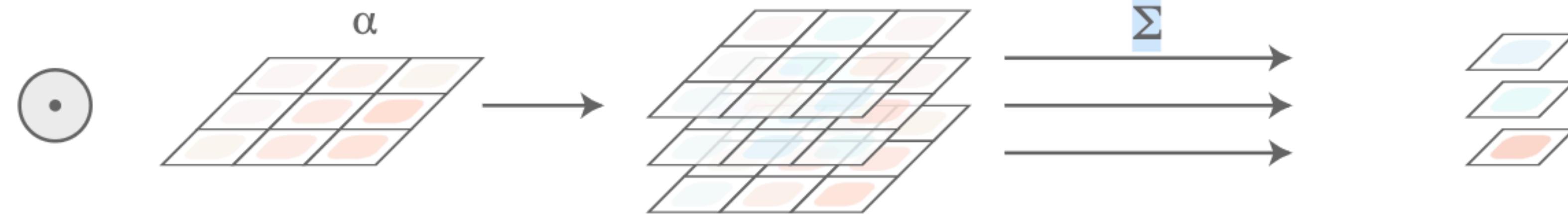




Attention vs FiLM



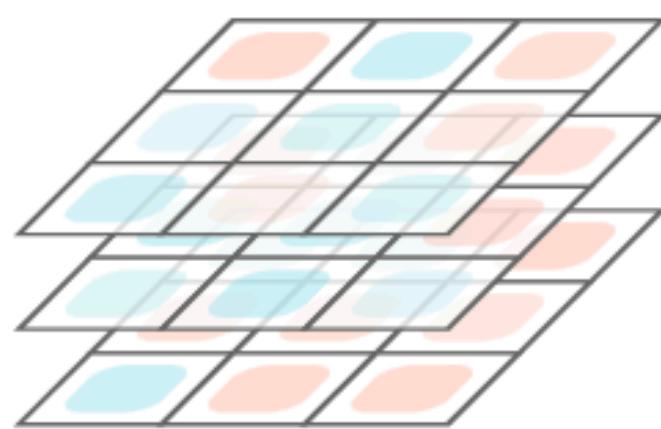
Attention computes a probability distribution over **locations**.



Attention pools over locations.

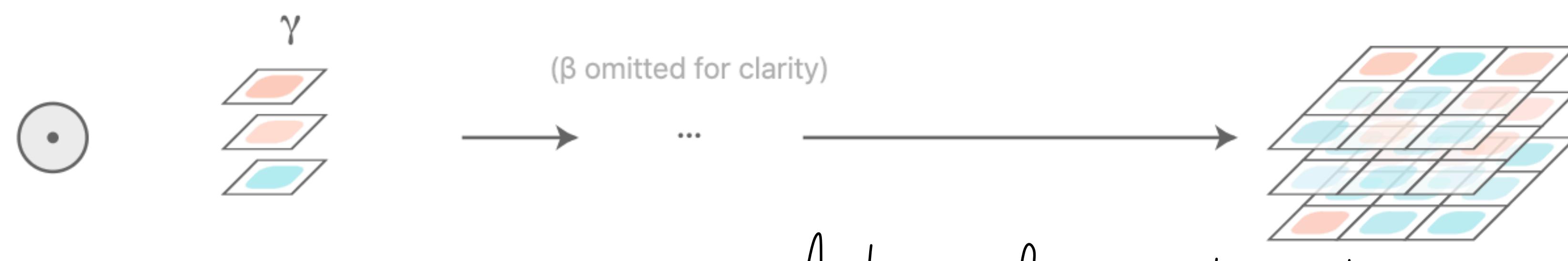
Attention summarizes the input into a vector.

„Specific locations / time-steps contain useful information“



FiLM computes a scaling vector applied to the **feature axis**.

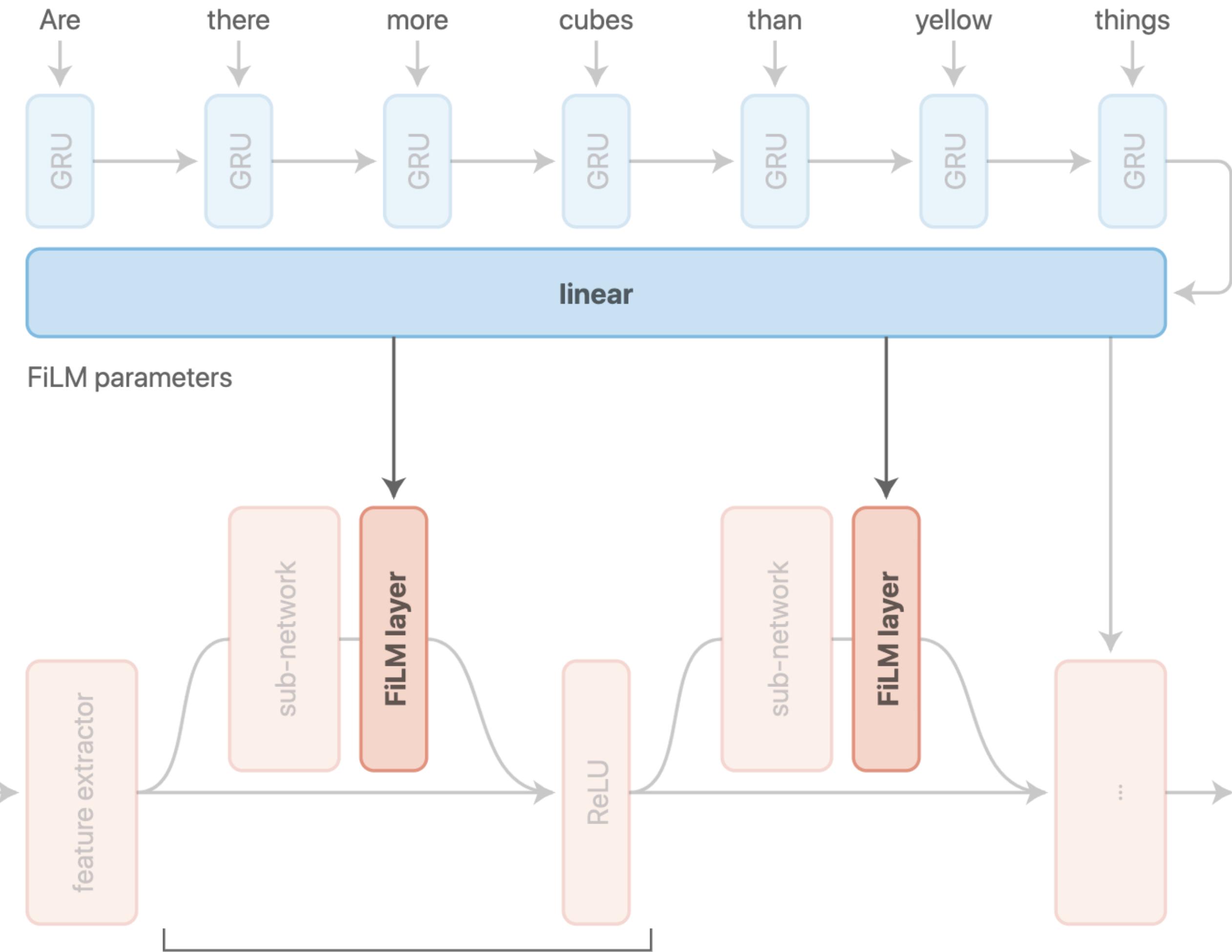
FiLM conserves input dimensions.



„Specific features contain useful information“

VQA

The **linguistic pipeline** acts as the FiLM generator.



FiLM layers in each residual block modulate the **visual pipeline**.

Learning visual reasoning without strong priors [\[PDF\]](#)

E. Perez, H. de Vries, F. Strub, V. Dumoulin, A. Courville.

ICML Workshop on Machine Learning in Speech and Language Processing. 2017.

FiLM: Visual Reasoning with a General Conditioning Layer [\[PDF\]](#)

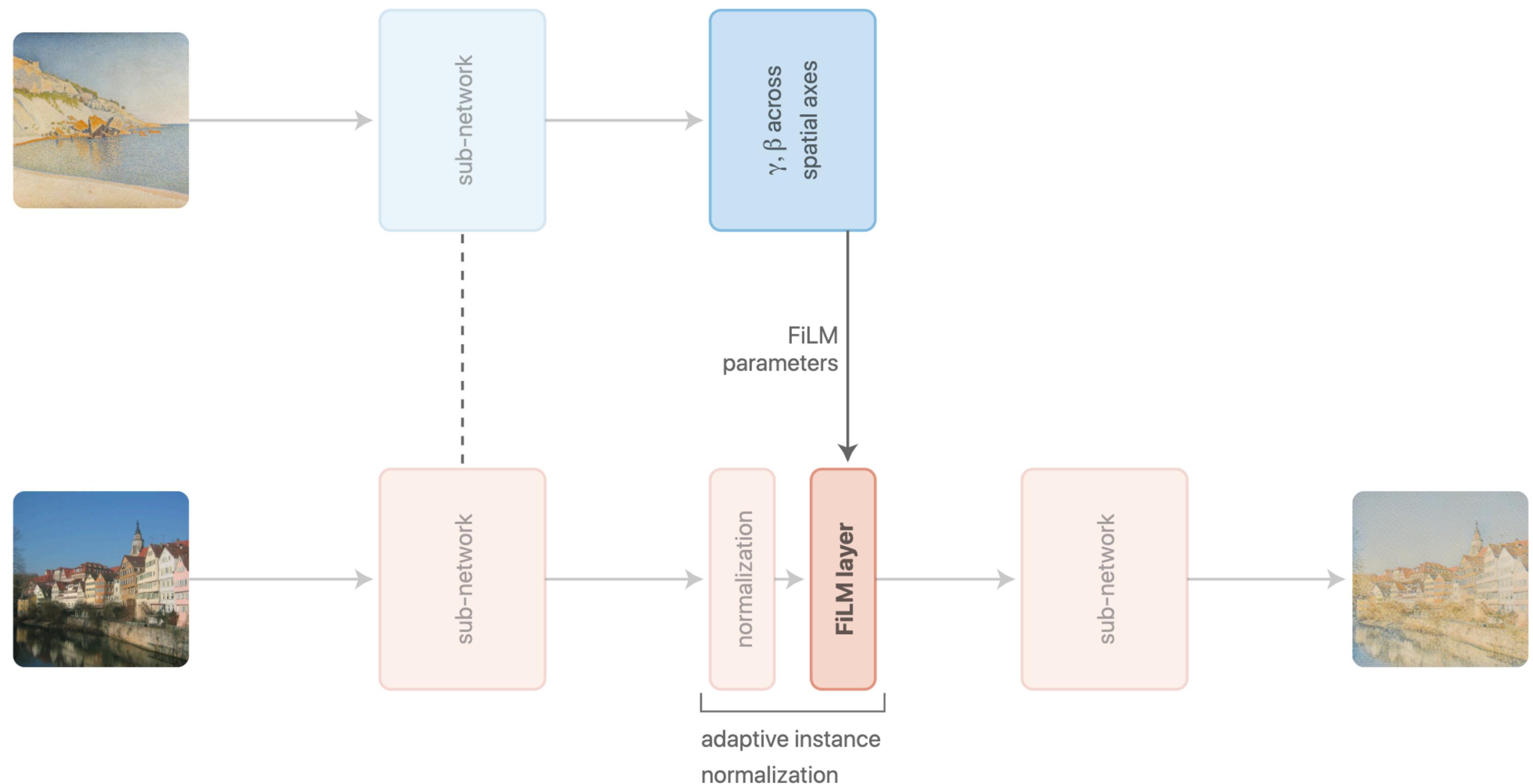
E. Perez, F. Strub, H.d. Vries, V. Dumoulin, A. Courville.

AAAI. 2018.

Style transfer

Conditional Batch
Normalization

$$z = \beta_s \left(\frac{x - \mu}{\sigma} \right) + \beta_s$$



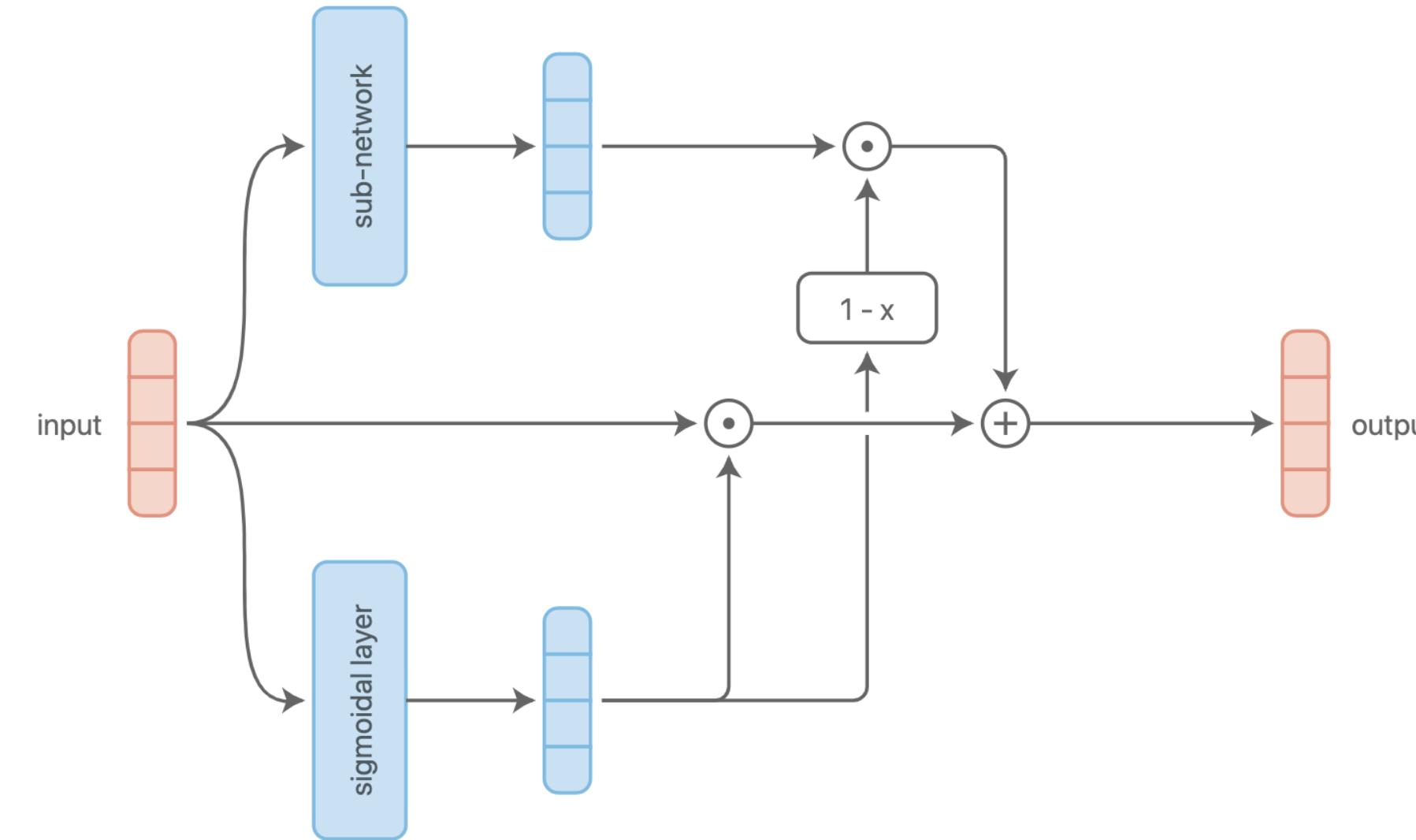
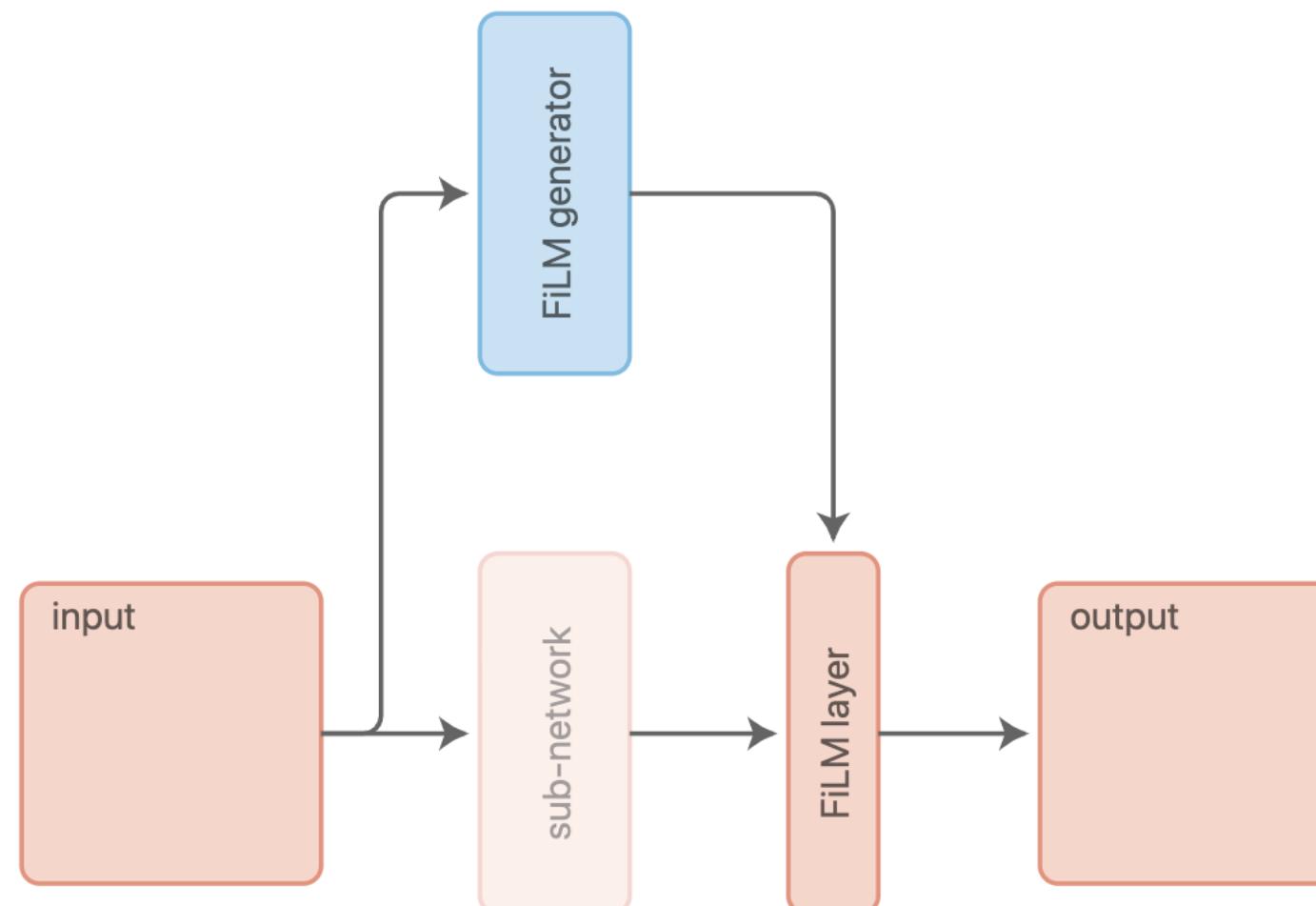
Arbitrary style transfer in real-time with adaptive instance normalization [\[PDF\]](#)

X. Huang, S. Belongie.

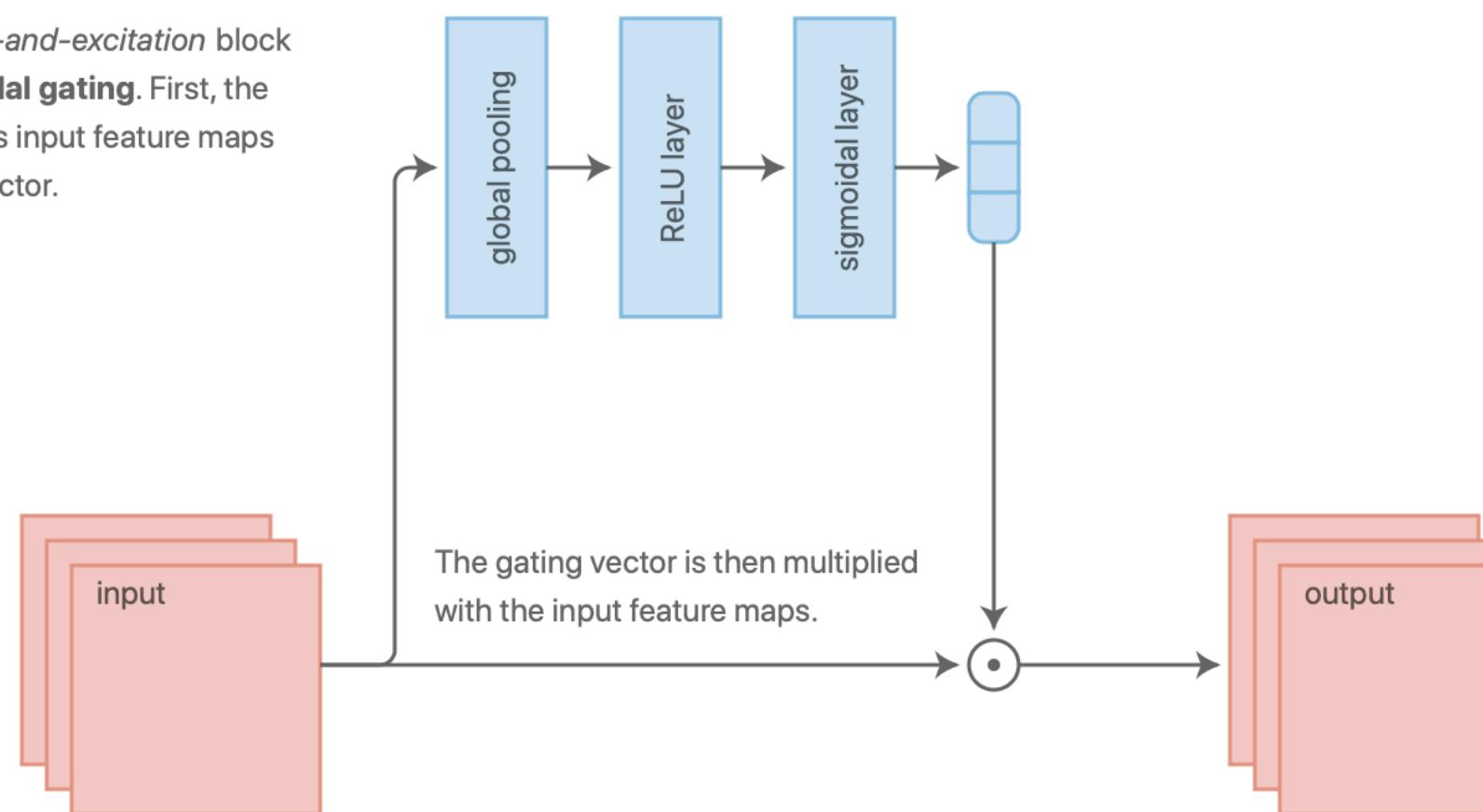
Proceedings of the International Conference on Computer Vision. 2017.

<https://distill.pub/2018/feature-wise-transformations/>

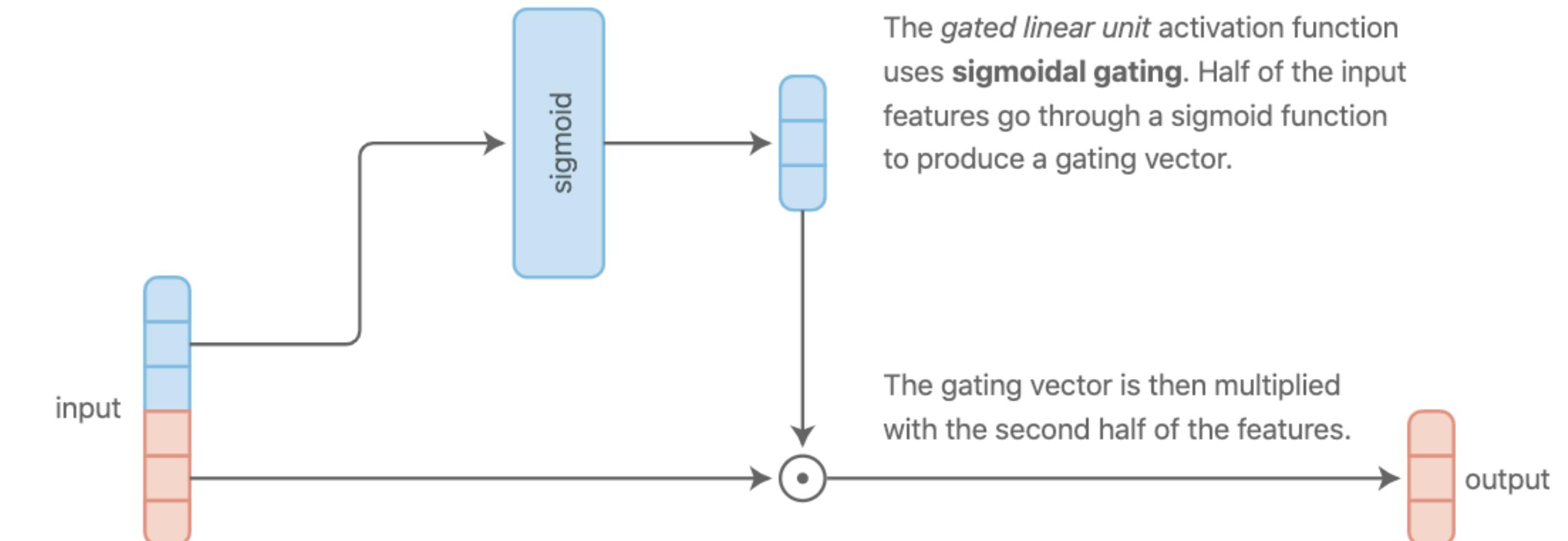
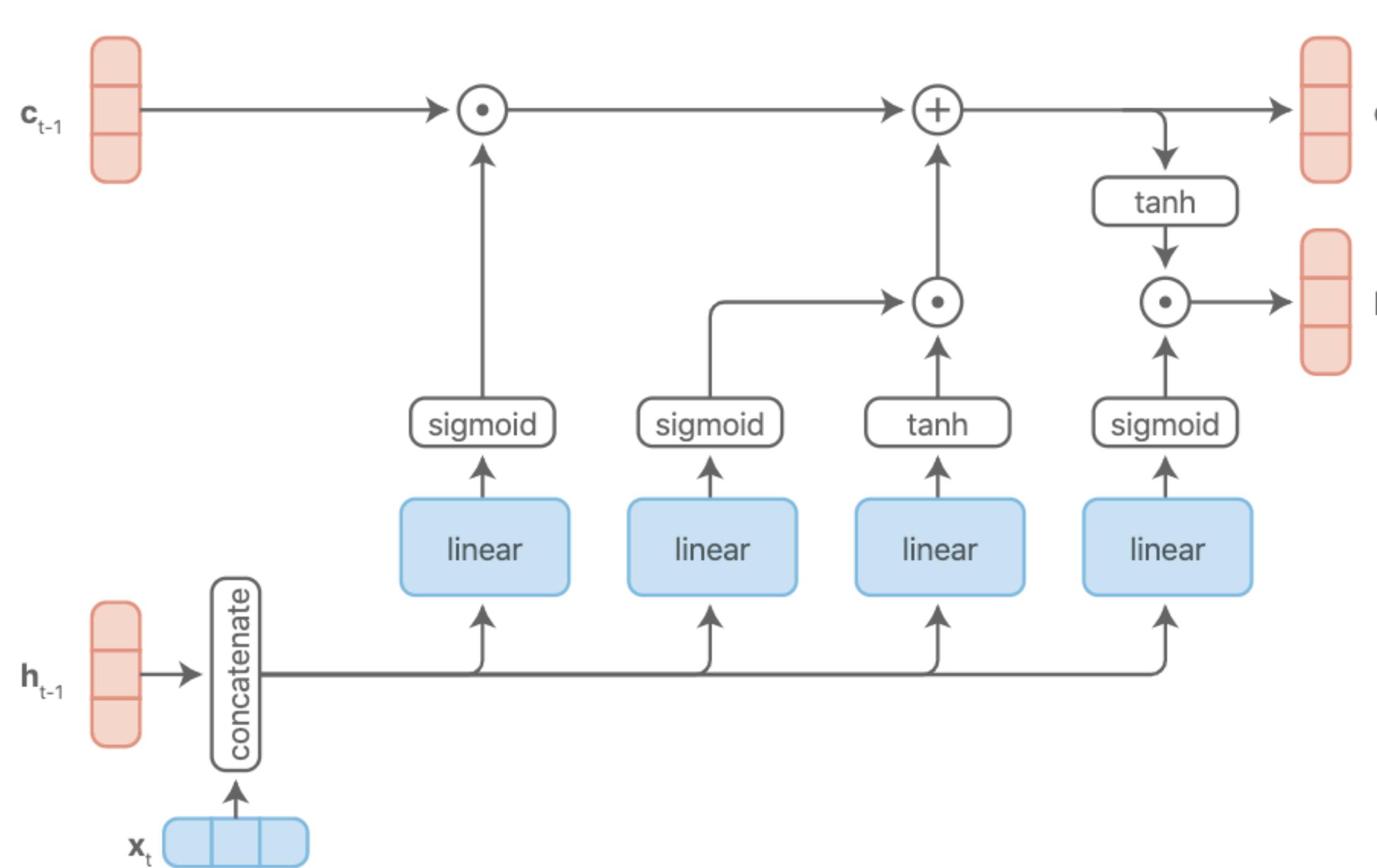
Image recognition: Highway Networks and Squeeze-and-excitation



The squeeze-and-excitation block uses **sigmoidal gating**. First, the network maps input feature maps to a gating vector.



NLP: LSTM, GLU and Gated Attention Reader

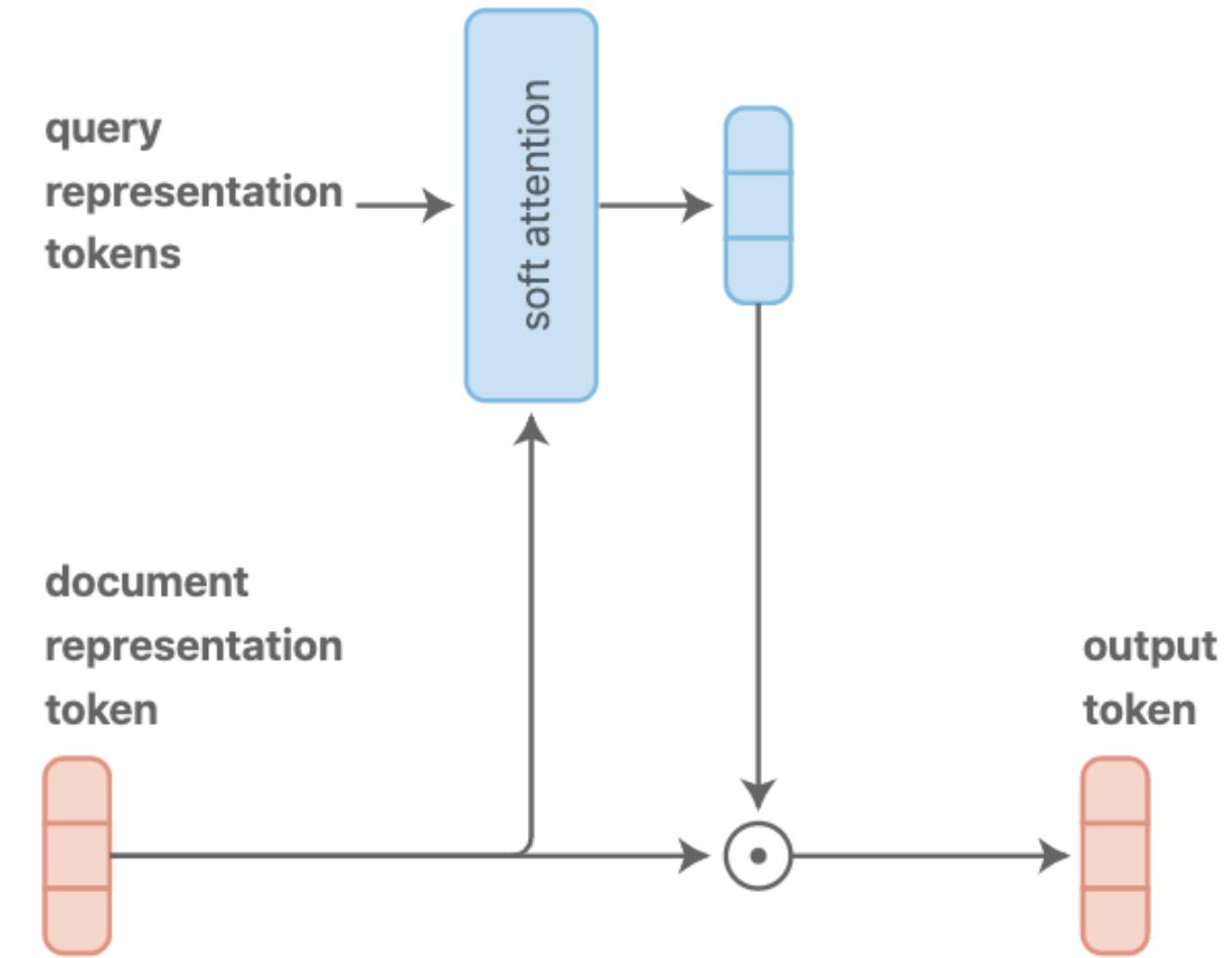


Dhingra et al. use **conditional scaling** to integrate query information into a document processing network. Applying soft attention to the **query representation tokens** produces the scaling vector.

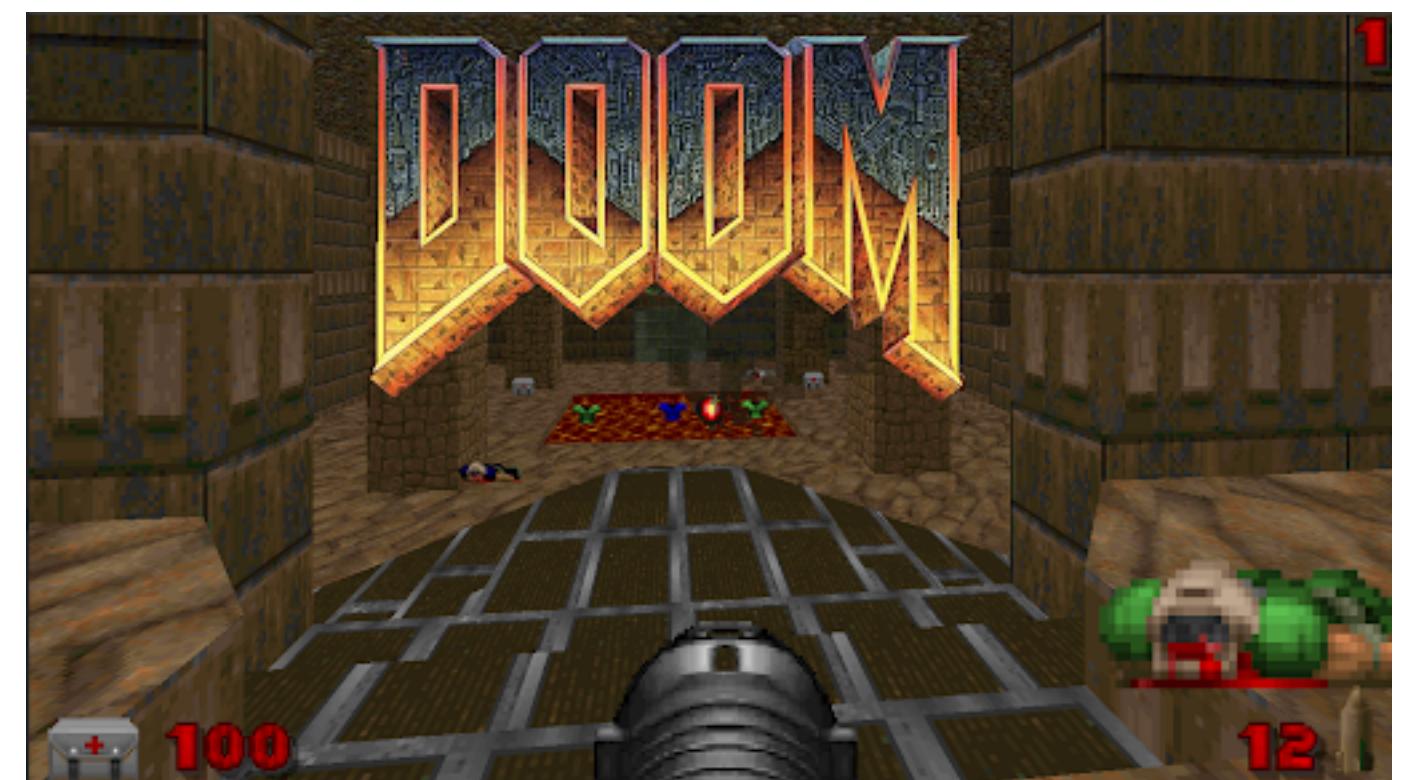
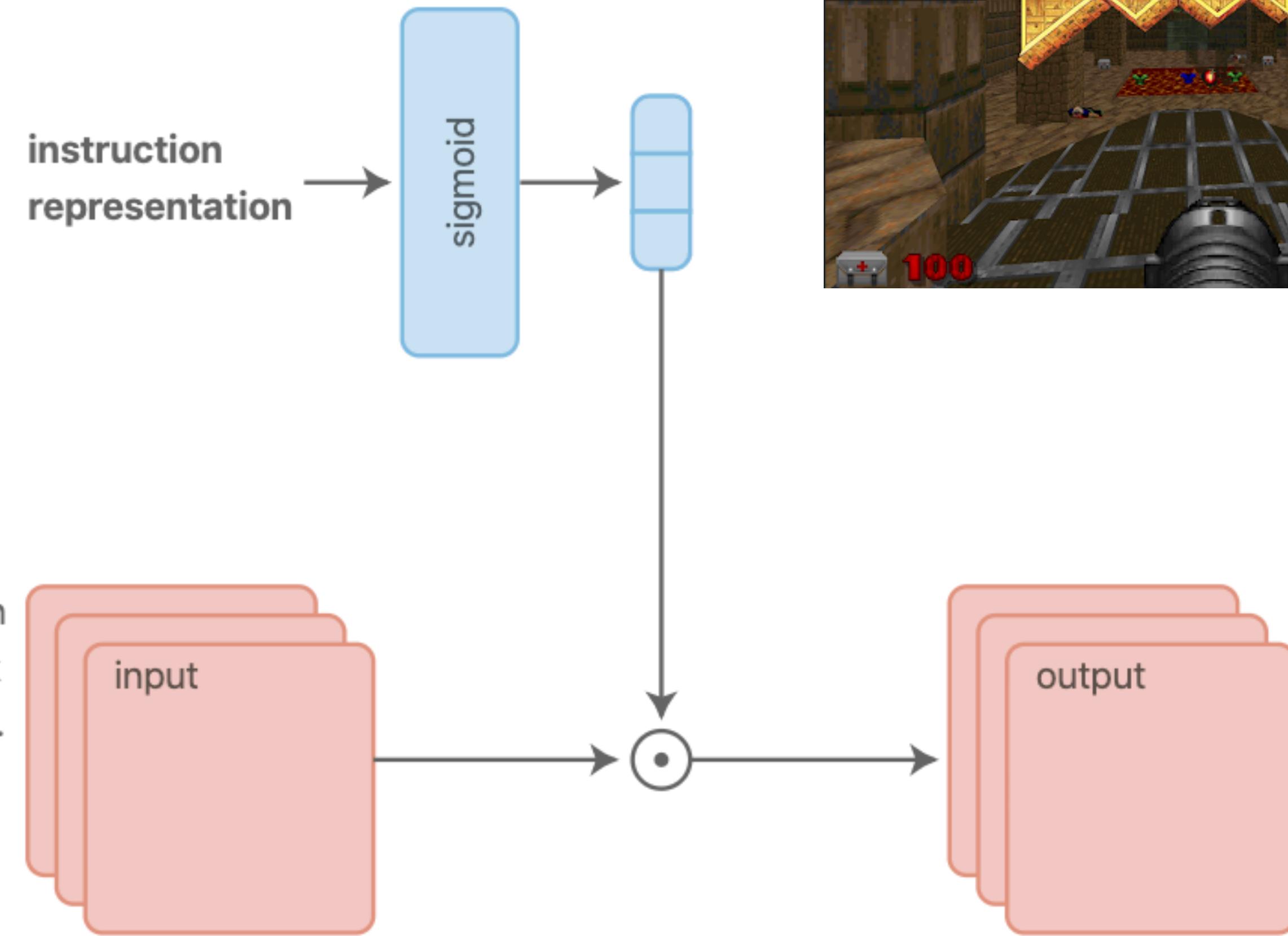
The scaling vector is then multiplied with the input **document representation token**

The *gated linear unit* activation function uses **sigmoidal gating**. Half of the input features go through a sigmoid function to produce a gating vector.

The gating vector is then multiplied with the second half of the features.



Chaplot et al. use **sigmoidal gating** as a multimodal fusion mechanism. An instruction representation is mapped to a scaling vector via a sigmoid layer.

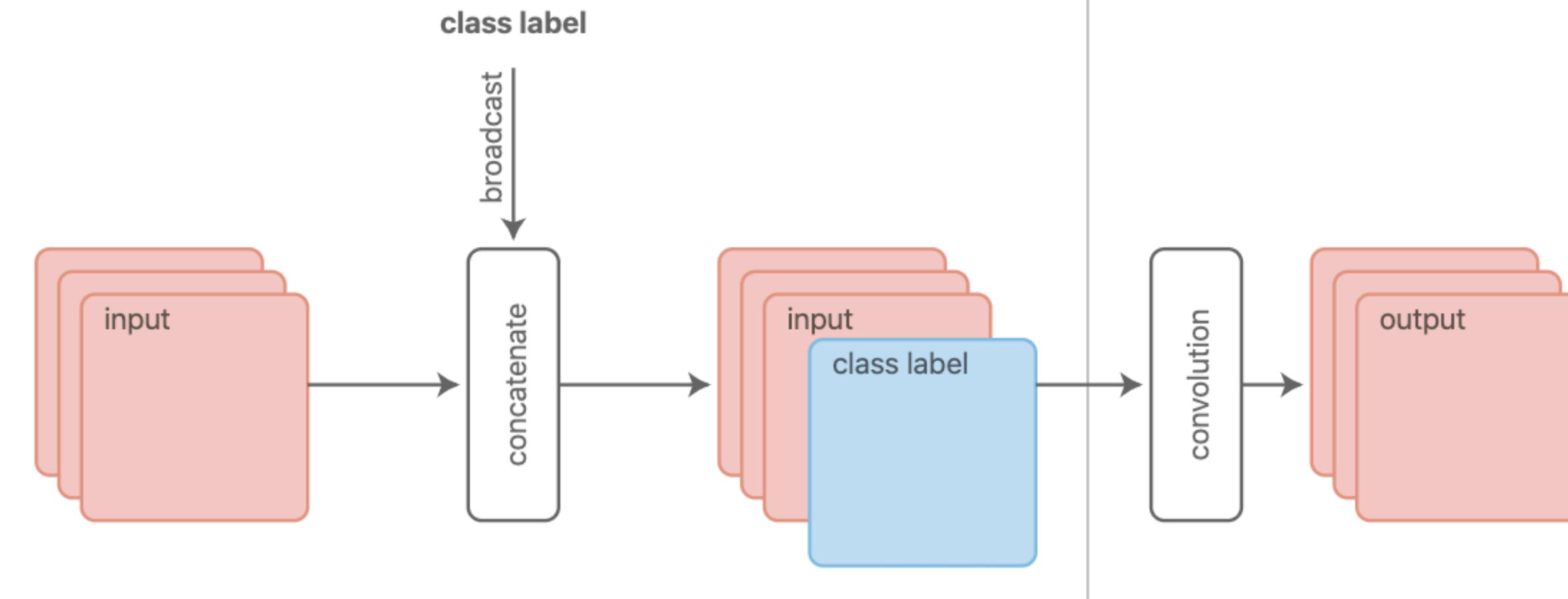


The scaling vector is then multiplied with the input feature maps. A policy network uses the result to decide the next action.

Generative modeling: DCGAN, PixelCNN

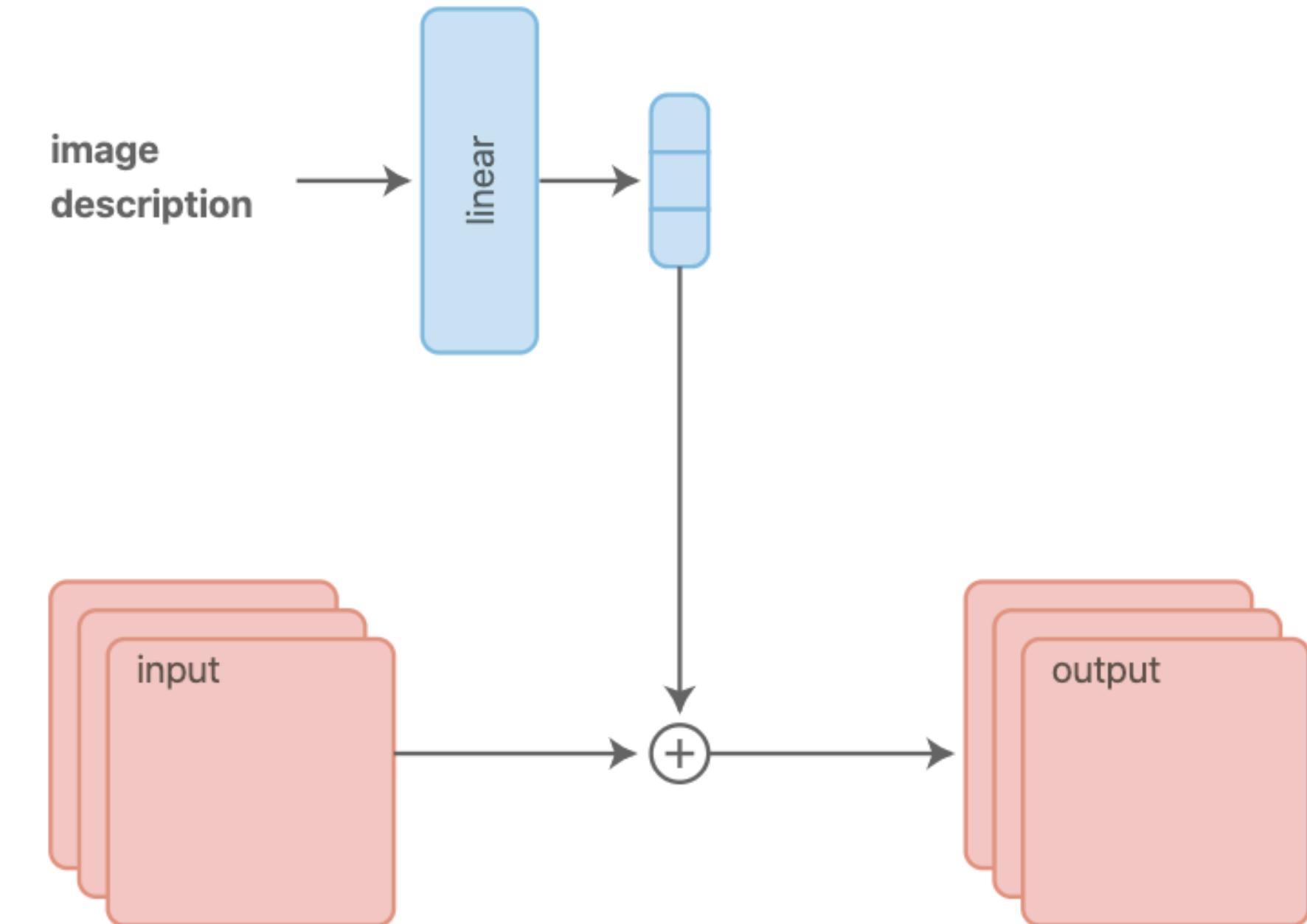
Concatenation-based conditioning is used in the class-conditional DCGAN model. Each convolutional layer is concatenated with the broadcasted label along the channel axis.

The resulting stack of feature maps is then **convolved** to produce the conditioned output.



PixelCNN uses **conditional biasing**.
The model first maps a high-level
image description to a bias vector.

Then, it adds the bias vector to
the input stack of feature maps
to condition convolutional layers.

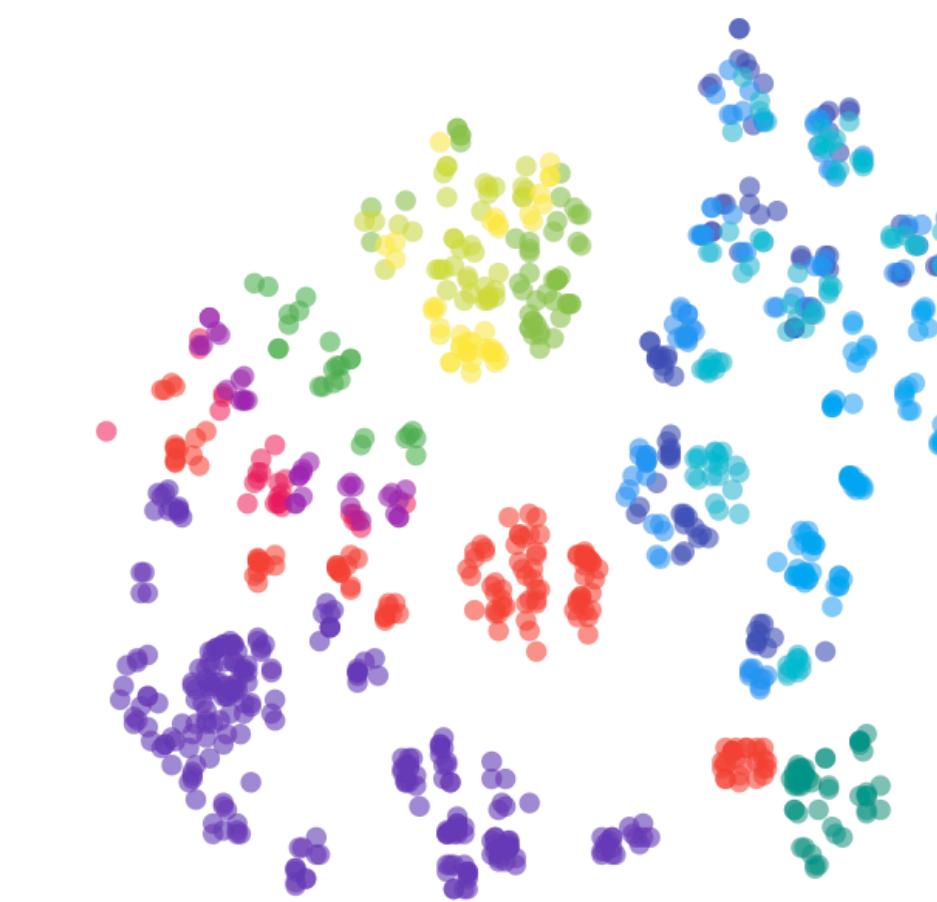


FiLM as a task representation

$$\text{FiLM}(\mathbf{x}) = \gamma(\mathbf{z}) \odot \mathbf{x} + \beta(\mathbf{z})$$

- each task
 - question about an image
 - a painting style to imitate
 - elicits a different set of FiLM parameters via the FiLM generator
 - representation in terms of how to modulate the FiLM-ed network

Visual reasoning model



Question type



- Exists
- Less than
- Greater than
- Count
- Query material
- Query size
- Query color
- Query shape
- Equal color
- Equal integer
- Equal shape
- Equal size
- Equal material

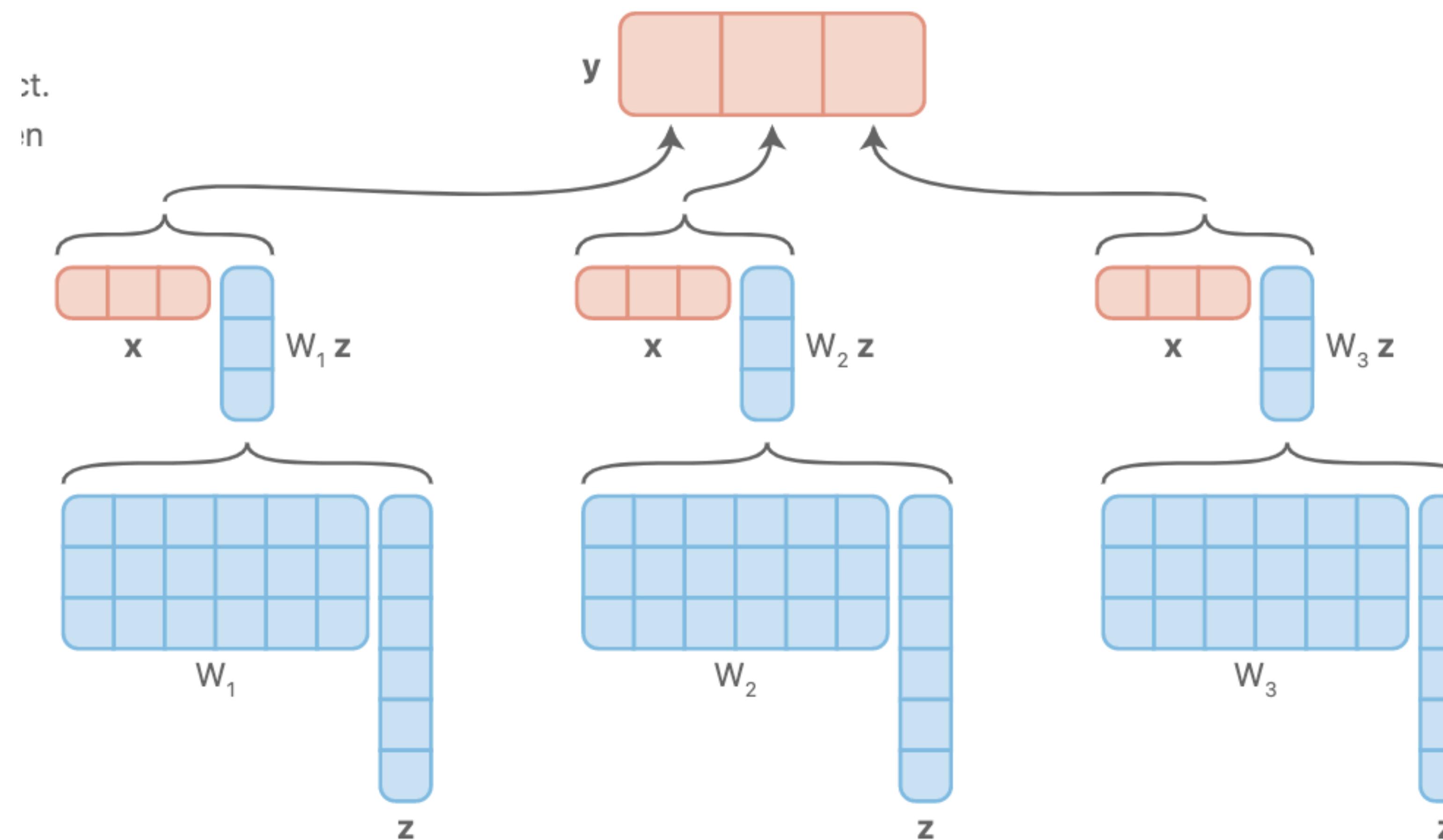
Reset pan / zoom

The background of the image is a vibrant, abstract swirl of liquid. It features a central cluster of bright red and orange hues, which transition into darker shades of red and then into a deep, rich blue on the right side. The liquid appears to be in motion, with visible streaks and ripples creating a sense of depth and fluidity.

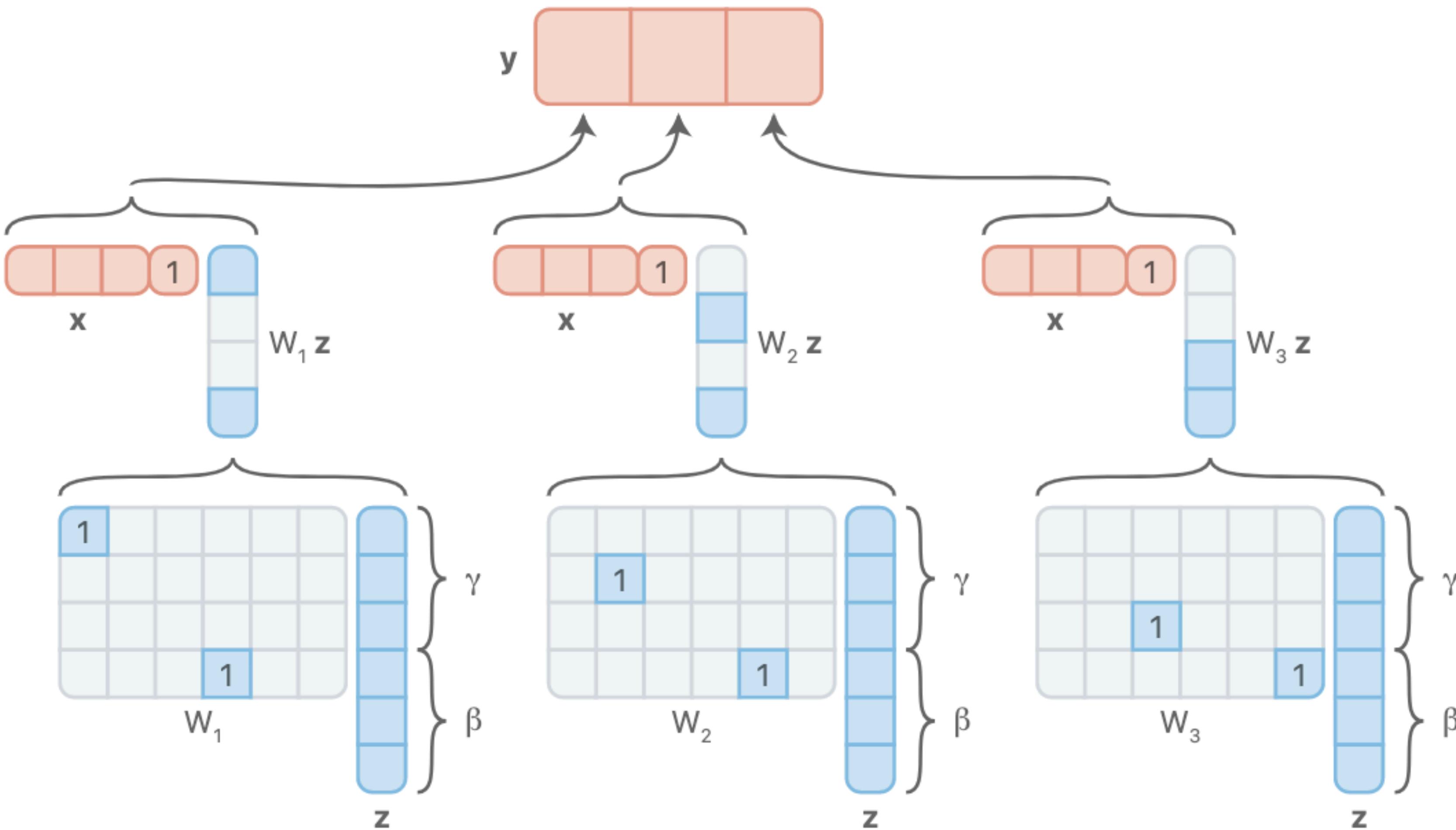
Billinear form

Bilinear form

$$y_k = \mathbf{x}^T W_k \mathbf{z}$$



FiLM: Low-Rank Bilinear form



Multiplicative Interactions

Additive form

Given arbitrary inputs \mathbf{x} and \mathbf{z} , where we are trying to model $f(\mathbf{x}, \mathbf{z})$, a single layer might typically look like:

$$f(\mathbf{x}, \mathbf{z}) = \mathbf{W}[\mathbf{x}; \mathbf{z}] + \mathbf{b}$$

where $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{z} \in \mathbb{R}^m$, $\mathbf{W} \in \mathbb{R}^{(m+n) \times k}$

Multiplicative Interactions

Instead, we now have the **bilinear** form:

$$f(\mathbf{x}, \mathbf{z}) = \mathbf{z}^T \mathbb{W} \mathbf{x} + \mathbf{z}^T \mathbf{U} + \mathbf{V} \mathbf{x} + \mathbf{b}$$

$\mathbf{x} \in \mathbb{R}^n$, $\mathbf{z} \in \mathbb{R}^m$, $\mathbb{W} \in \mathbb{R}^{m \times n \times k}$ and $(\mathbf{z}^T \mathbb{W} \mathbf{x})_k = \sum_{ij} \mathbf{z}_i \mathbb{W}_{ijk} \mathbf{x}_j$.

HyperNetworks

$$f(\mathbf{x}, \mathbf{z}) = \mathbf{z}^T \mathbb{W} \mathbf{x} + \mathbf{z}^T \mathbf{U} + \mathbf{V} \mathbf{x} + \mathbf{b}$$

Set $\mathbf{W}' = \mathbf{z}^T \mathbb{W} + \mathbf{V}$ and $\mathbf{b}' = \mathbf{z}^T \mathbf{U} + \mathbf{b}$

$$f(\mathbf{x}, \mathbf{z}) = \mathbf{W}' \mathbf{x} + \mathbf{b}'$$

Implementation

```
# Simple python code for MI Layers
import sonnet as snt
import tensorflow as tf

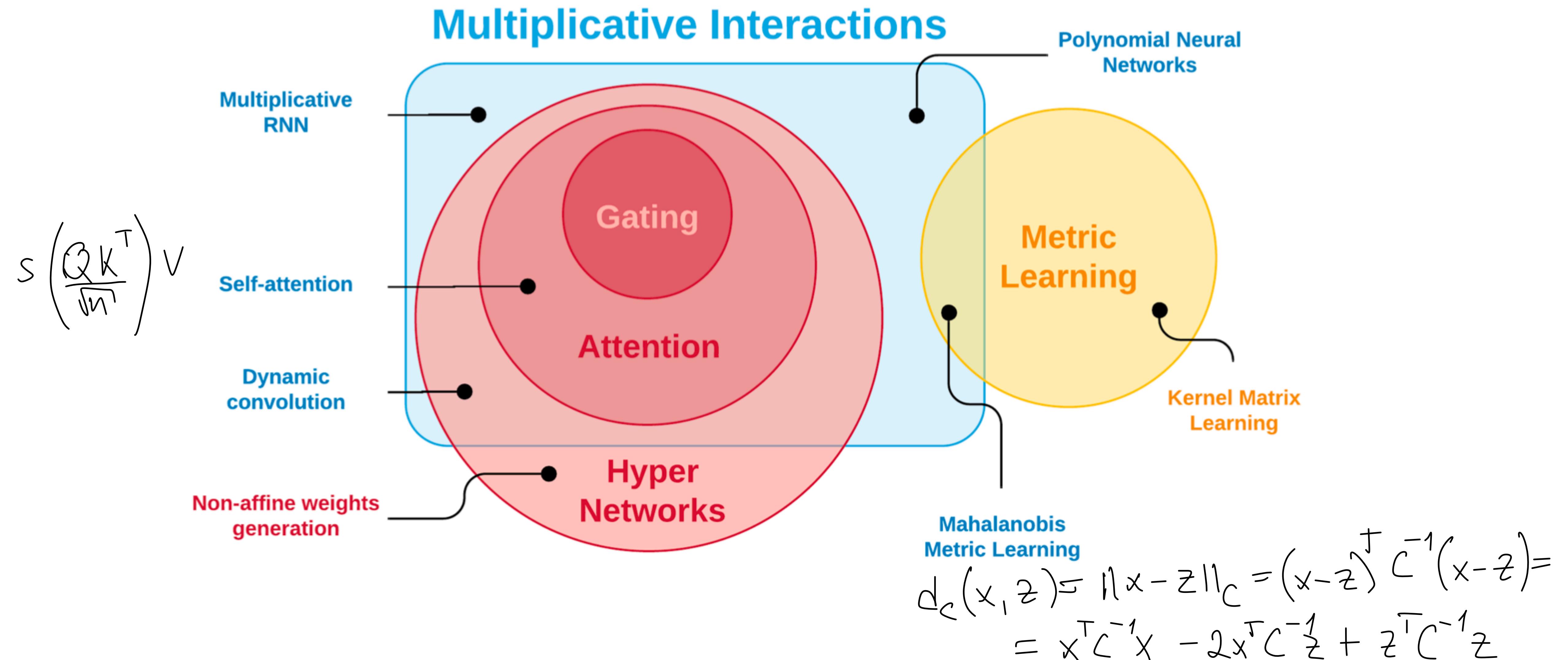
# A standard linear layer
# B is the batch size
# E is the input size
# C is the context size
x = ... # input of size [B, E]
z = ... # context of size [B, C]
xz = tf.concat([x,z], 1)
y = snt.Linear(output_size)(xz)

# Instead , we generate a W and b
# This defines an implicit 3D weight tensor
W = snt.Linear(output_size * input_size)(z)
b = snt.Linear(output_size)(z)

# Reshape to the correct shape
# Note: we have B weight matrices
# i.e. one per batch element
W = tf.reshape(W, [B, input_size, output_size])

# Output
y = tf.matmul(x, W) + b
```

Multiplicative Interactions



Order of multiplicative interactions

	Wz^T		
xWz^T			
Multiplicative interaction	Scaling	Hadamard product	General bilinear form
Projection matrix class	Scalar	Diagonal	Unconstrained
HyperNetwork output	Scalar	Vector	Matrix

Diagram illustrating the order of multiplicative interactions for the expression xWz^T across three stages:

- Stage 1:** Wz^T is shown as a 2x2 matrix multiplied by a scalar (blue square) and a 2x2 diagonal matrix (gray with 1s).
- Stage 2:** Wz^T is shown as a 2x2 matrix multiplied by a vector (blue vertical stack) and a 2x2 diagonal matrix (gray with 1s).
- Stage 3:** Wz^T is shown as a 2x2 matrix multiplied by a 2x2 matrix (blue and gray).

Below the stages, the expression xWz^T is expanded into three components:

- Scaling:** A green horizontal vector (green blocks) multiplied by a 2x2 matrix where the top-left element is orange and the other elements are white.
- Hadamard product:** A green horizontal vector (green blocks) multiplied by a 2x2 matrix where the top-left element is orange, the top-right is red, the bottom-left is red, and the bottom-right is pink.
- General bilinear form:** A green horizontal vector (green blocks) multiplied by a 2x2 matrix where all elements are pink.

Dashed boxes group the components corresponding to each stage.

Expressivity of the models

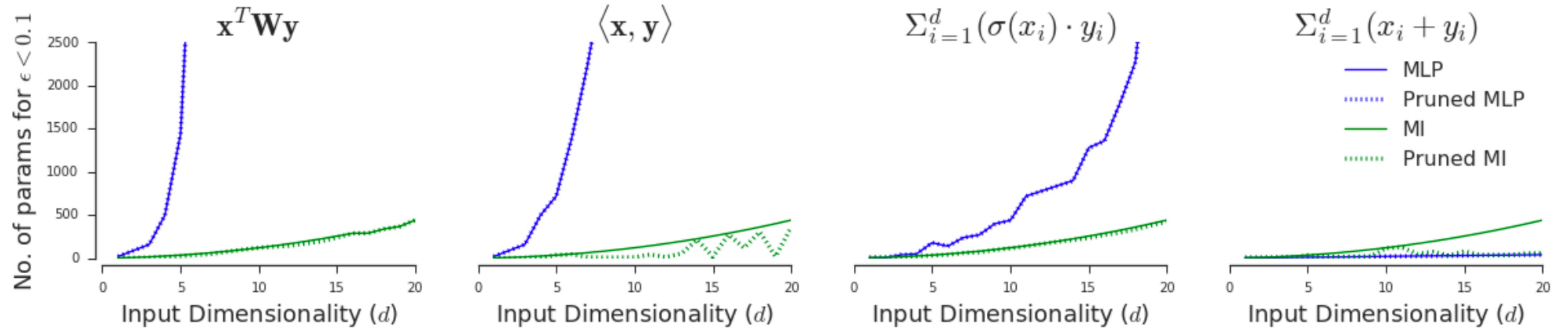


Figure 2: Number of parameters needed for a regular, single layer MLP (blue line) to represent the function up to 0.1 MSE over the domain of a standard d -dimensional Gaussian compared to the same quantity for a multiplicative model (green line). σ denotes sigmoid. Dotted lines represent pruned models where all weights below absolute value of 0.001 were dropped. Note that for MLP all parameters are actually used, while for MI module some of these functions (summation and dot product) can be compactly represented with pruning.

Toy Multi-task Learning

$$y = a_i x + b_i \text{ and } y = a_i \sin(10x) + b_i$$

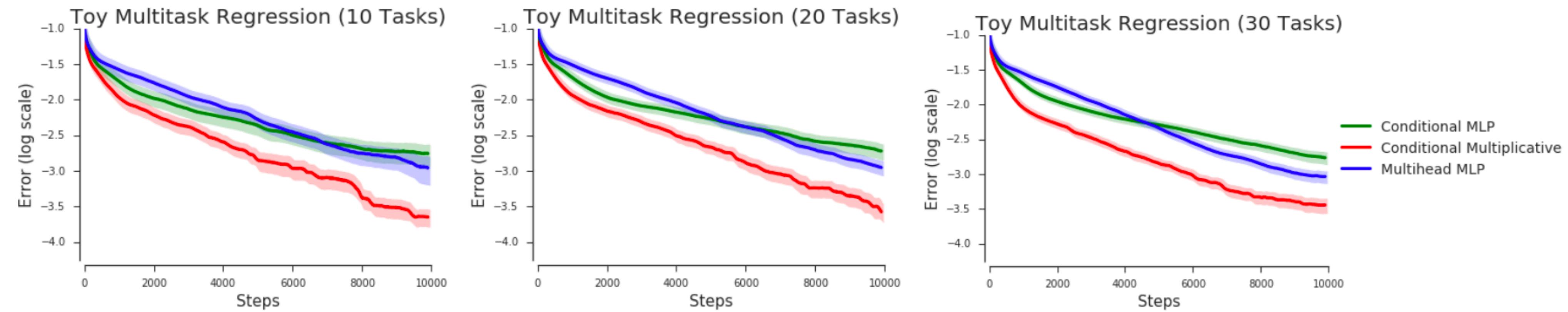


Figure 3: Averaged learning curves for different models while varying the number of tasks in the toy multitask regression domain. Shaded regions represent standard error of mean estimation.

Multi-Task RL on DMLab-30

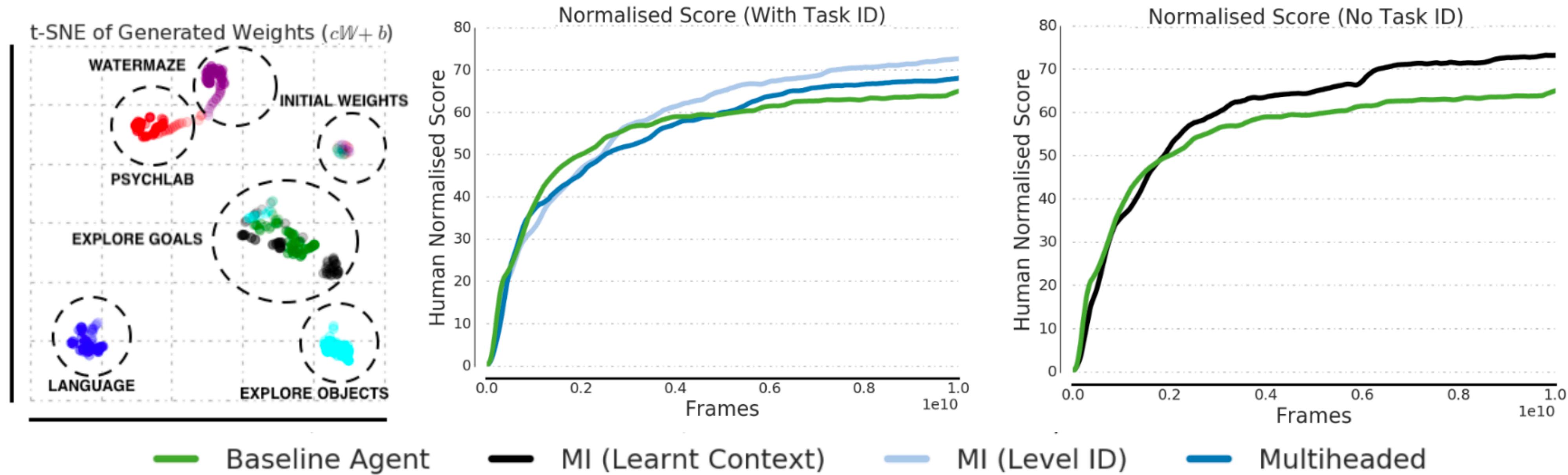


Figure 4: (a) A t-SNE plot of the generated weights from an \mathcal{M} layer. (b) Human normalised performance (capped at 100) when using task ID as context to an \mathcal{M} layer. (c) Using a learnt context instead.

Language Modeling

Table 1: Word-level perplexity on WikiText-103

	Model	Valid	Test	No. Params
LSTM	Rae et al. (2018)	34.1	34.3	88M
Gated CNN	Dauphin et al. (2017)	-	37.2	-
RMC	Santoro et al. (2018)	30.8	31.6	-
Trellis Networks	Bai et al. (2019)	-	30.35	180M
TransformerXL	Dai et al. (2018)	17.7	18.3	257M
LSTM (ours)				
LSTM + MultDec				
LSTM + MultEncDec				

The background of the image is a vibrant, abstract pattern of ink droplets suspended in water against a black background. The ink forms intricate, flowing shapes in shades of blue, red, orange, and yellow. A large, semi-transparent blue sphere is positioned in the upper left quadrant. In the lower right quadrant, there is a dense, textured cluster of orange and yellow ink. The overall effect is organic and fluid, resembling a microscopic view of a biological specimen or a celestial body.

Thank you