

# GATO

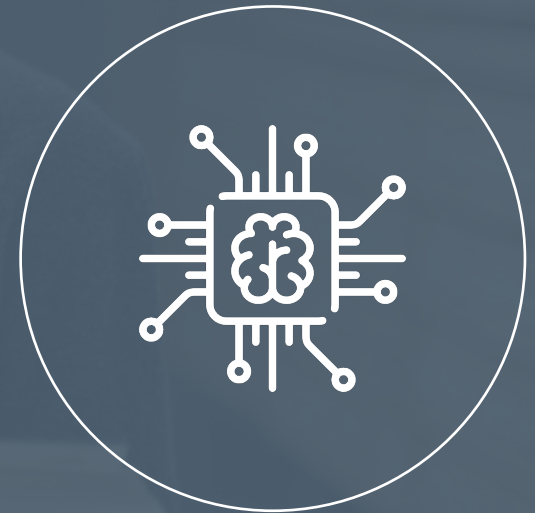
## A Scalable Real-time Architecture for Learning Knowledge from Unsupervised Sensorimotor Interaction

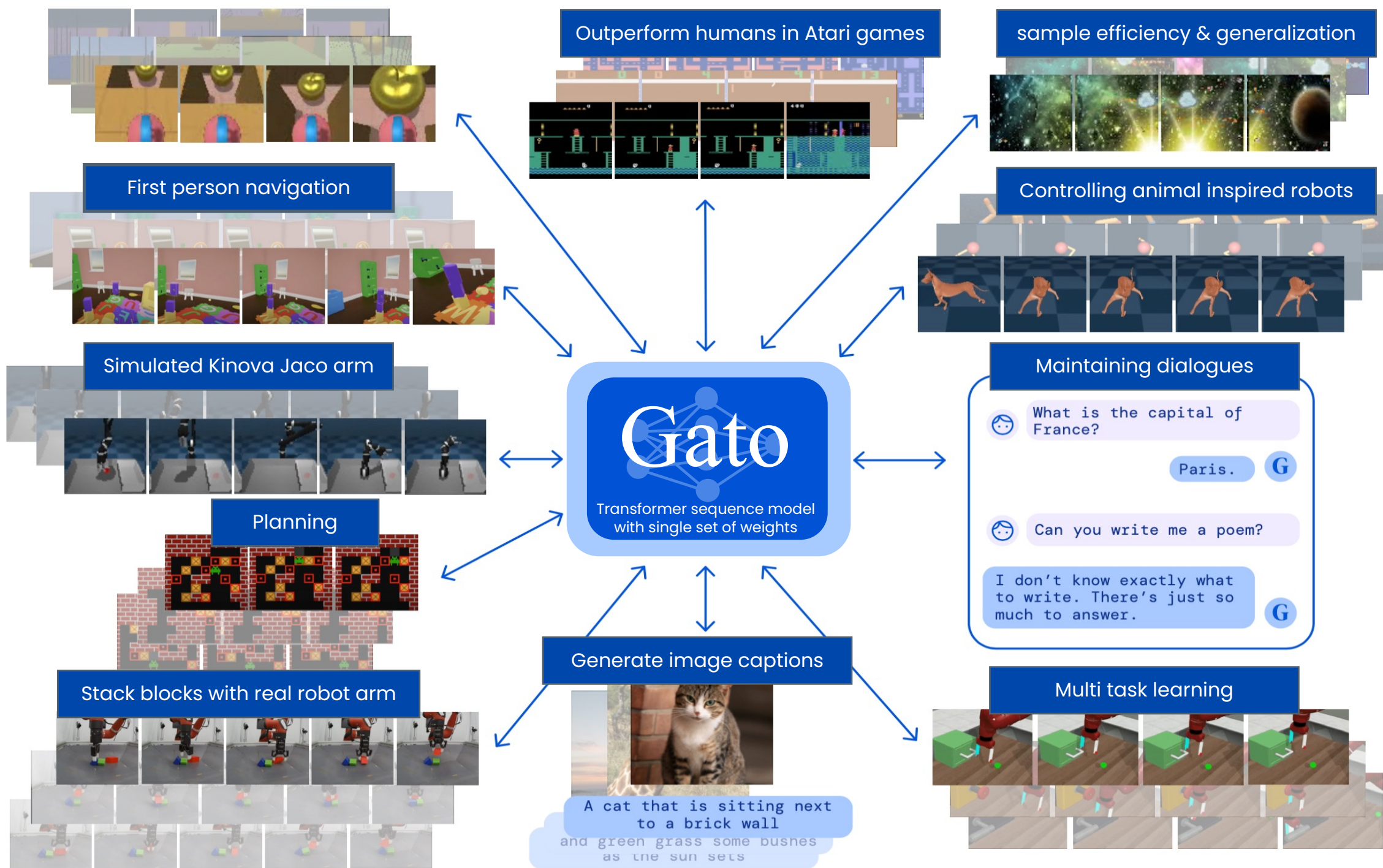
---

Token embedding and sequencing



**Niklas Forsstroem**  
forsstroemniklas@gmail.com

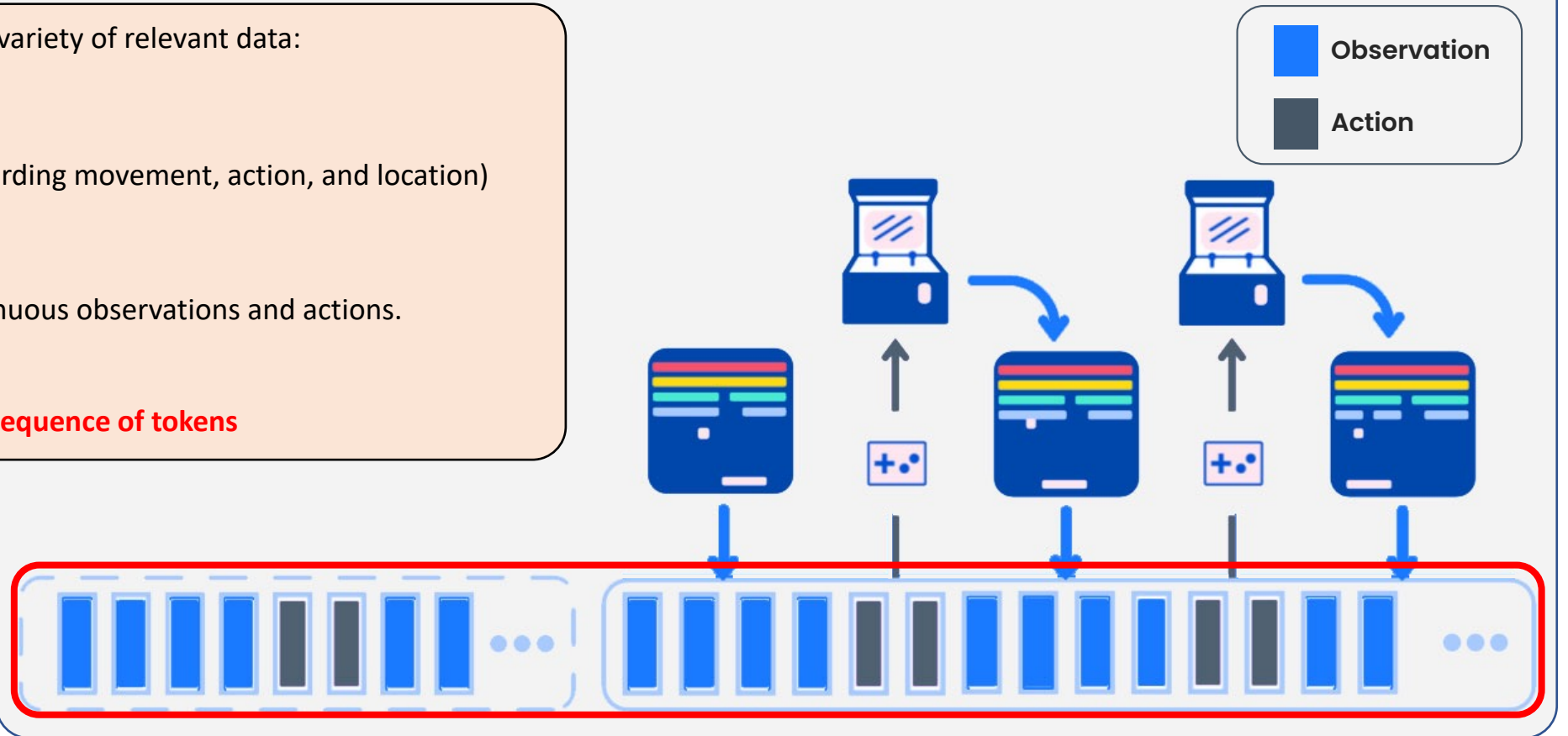




# Model application

## General approach

- Trained on the widest possible variety of relevant data:
  - Images
  - Text
  - Proprioception (info regarding movement, action, and location)
  - Joint torques
  - Button presses
  - Other discrete and continuous observations and actions.
- All data is serialized into a flat **sequence of tokens**





# High-level tokenization

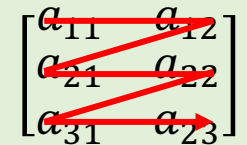
## Sequencing of agent data

- **Episodes** are presented to the agent in order of time.
- **Each timestep** is presented in the order: **Observations** ( $[y_{1:k}, x_{1:m}, z_{1:n}]$ ), **Separator** ('|'), **Actions** ( $a_{1:A}$ )
- A sequence of tokens is the concatenation of data from T timesteps:

$$s_{1:L} = \left[ [y_{1:k}^1, x_{1:m}^1, z_{1:n}^1, '|', a_{1:A}^1], \dots, [y_{1:k}^T, x_{1:m}^T, z_{1:n}^T, '|', a_{1:A}^T] \right]$$

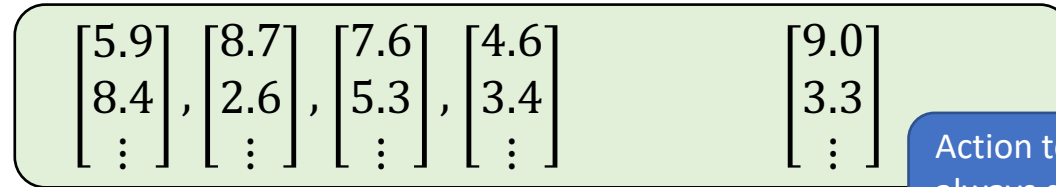
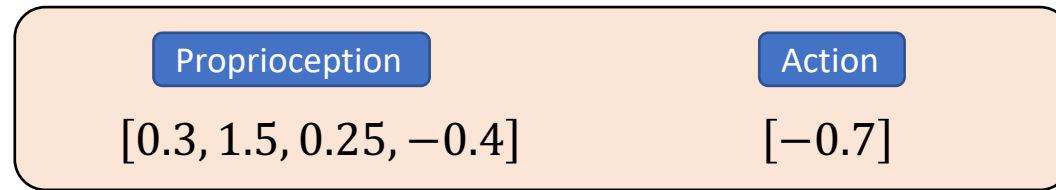
- **Observations** ( $[y_{1:k}, x_{1:m}, z_{1:n}]$ ) are ordered according to a key, each item is sequenced as follows:
  - **Text tokens** ( $y_{1:k}$ ) are in the same order as the raw input text (with SentencePiece encoding).
  - **Image patch tokens** ( $x_{1:m}$ ) are in **row-major order**.
  - **Tensors** ( $z_{1:n}$ ) (such as discrete and continuous observations) are in **row-major order**.
- **Separator** ('|'); a designated separator token is provided after observations.
- **Actions** ( $a_{1:A}$ ) are tokenized as discrete or continuous values and in **row-major order**.

Row-major / raster order



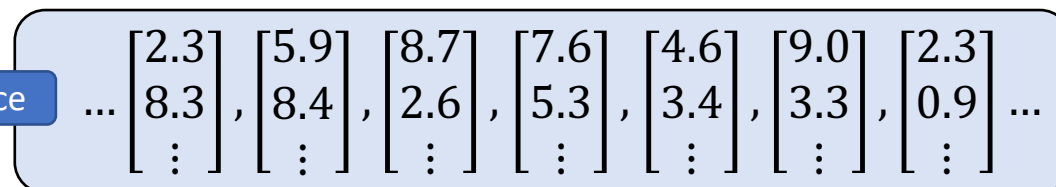
# Summarized state representations

## Tokenization and embedding



Action tokens  
always get same  
position encoding

Full sequence



Mu-law encoding & vector embedding

$$f(1.5) =$$

Uses lookup table

Local position encoding

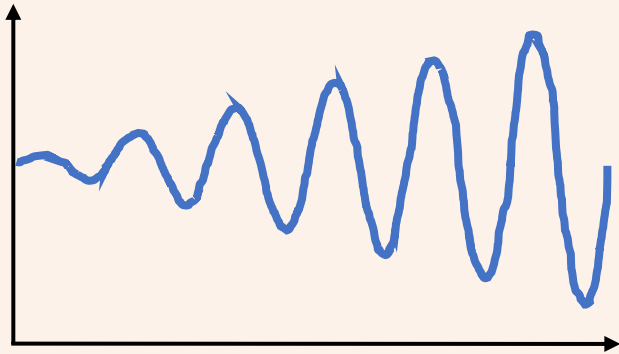
$$h\left(\begin{bmatrix} 5.9 \\ 8.4 \\ \vdots \end{bmatrix}, \begin{bmatrix} 8.7 \\ 2.6 \\ \vdots \end{bmatrix}\right) = \begin{bmatrix} 5.9 \\ 8.4 \\ \vdots \end{bmatrix} + \mathbf{p}_0, \quad \begin{bmatrix} 8.7 \\ 2.6 \\ \vdots \end{bmatrix} + \mathbf{p}_1$$

e.g.  $\mathbf{p}_0 = \begin{bmatrix} 3.0 \\ 5.8 \\ \vdots \end{bmatrix}$

# Mu-law companding

## Reducing dynamic range

Input signal



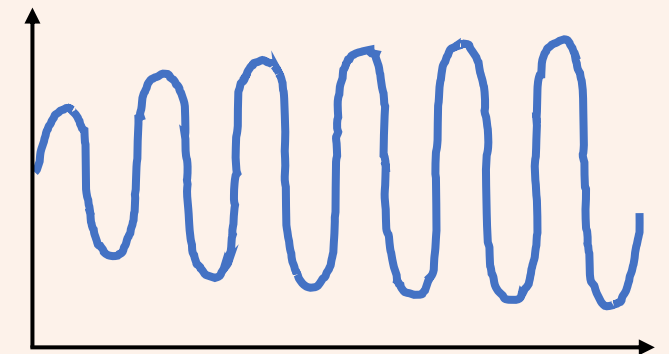
Preserves sign

$$F(x) = \text{sgn}(x) \frac{\log(|x|^\mu + 1.0)}{\log(M^\mu + 1.0)}$$

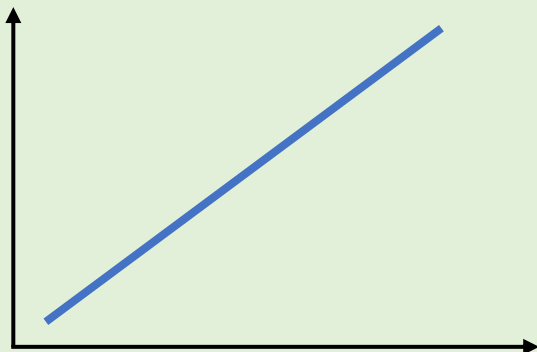
Constant denominator

$\mu = 100$   
 $M = 256$

Reduced dynamic range

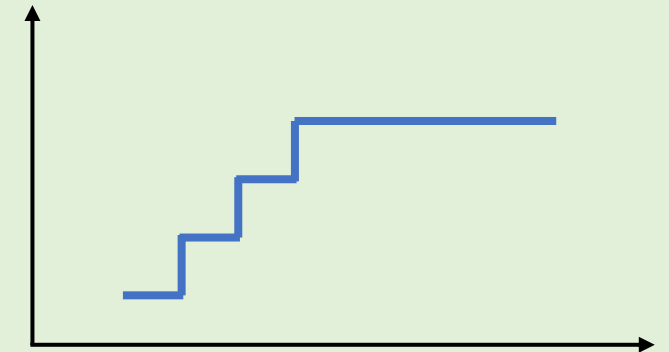


After mu-law companding

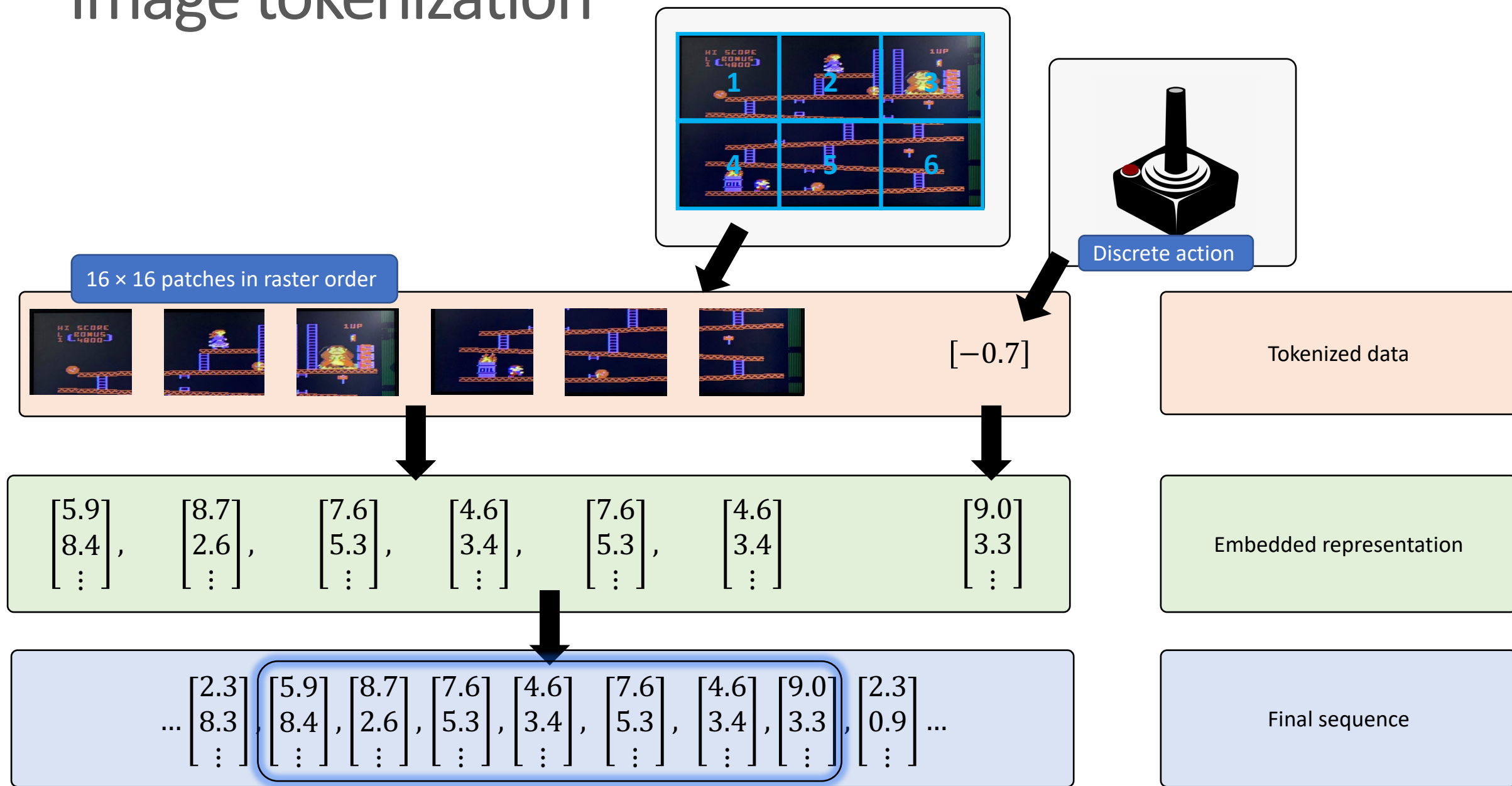


- Capping to  $[-1, 1]$
- Discretization into 1024 bins

Processed signal

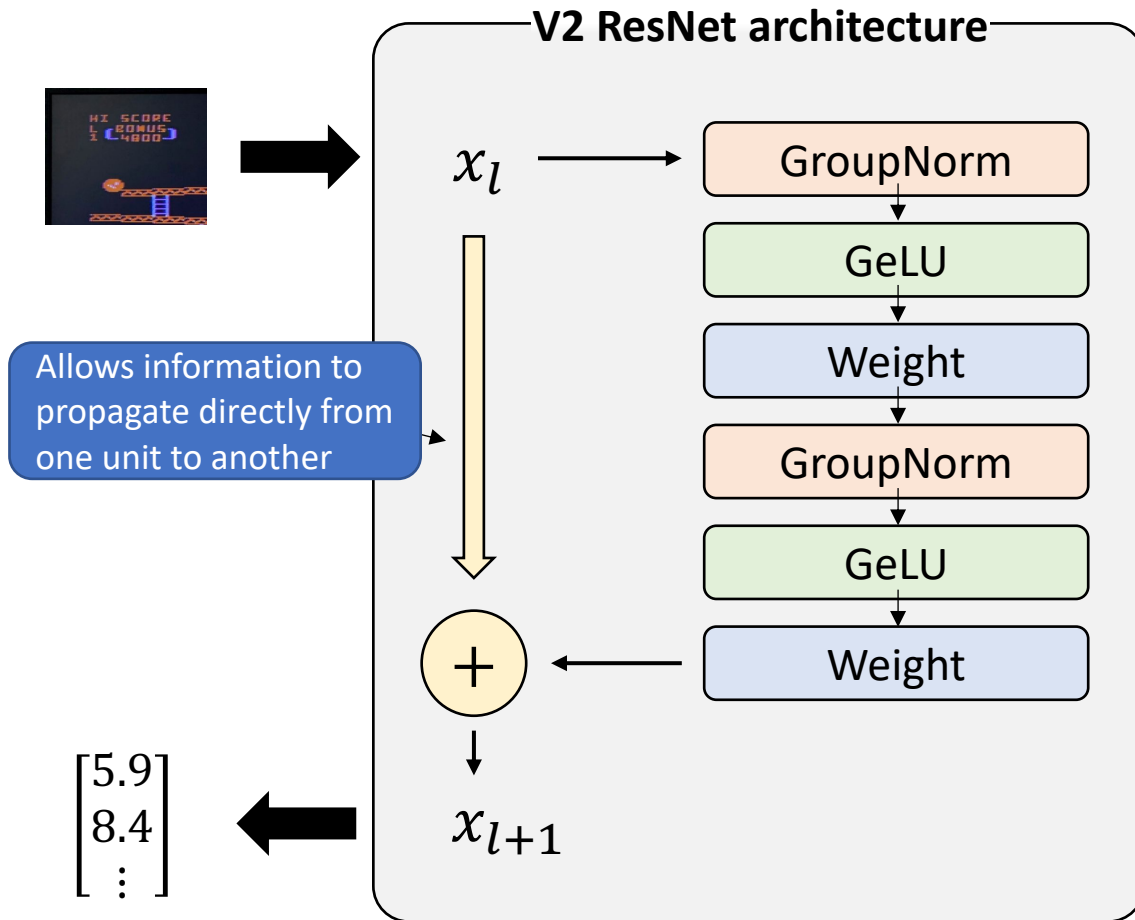


# Image tokenization



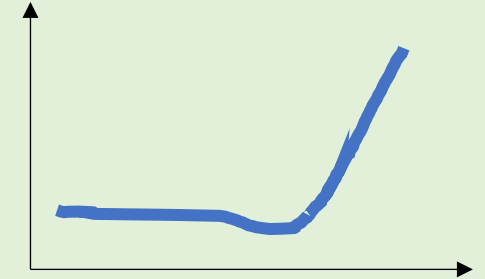
# Converting tokenized image to embedding

## ResNet architecture



- **GroupNorm** normalizes data to appropriate scales, similar to LayerNorm and InstanceNorm
- Good for distributed systems

- **GeLU** is an activation function.
- Provides well defined gradient in negative regime
- Prevents neurons from dying while bounding influence of negative regime activations



- The **weight** layers are configurable with hyper parameters
- Common to use small convolutional layers
- Multiple layers of these units can be stacked to increase complexity

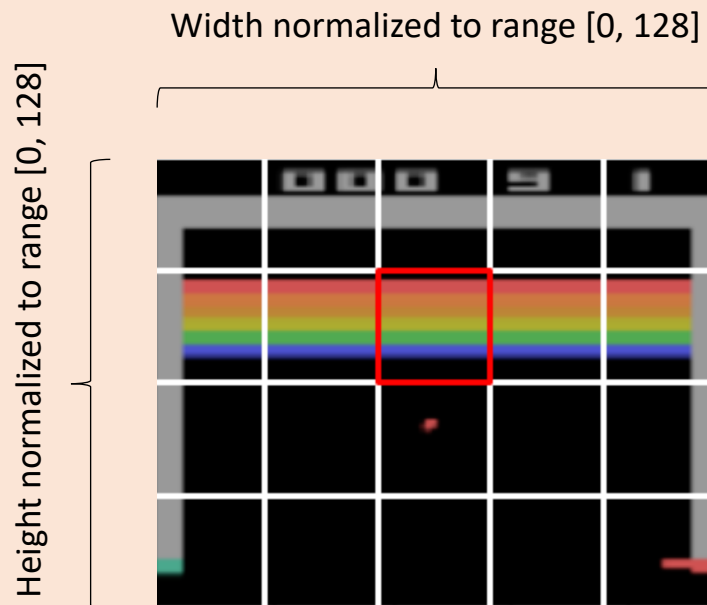


## 2.2 Embedding input tokens and setting output targets

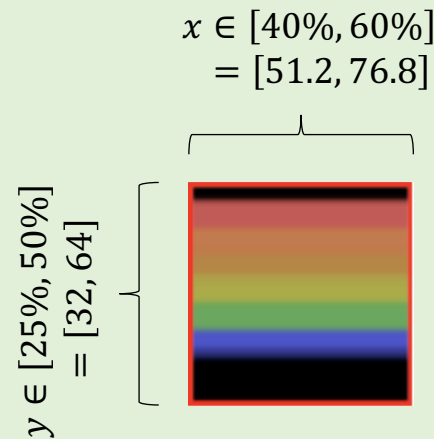
### Patch position encodings: Input from image patches

- For image patch token embeddings, between the embedded image and mapping onto the final sequence, a learnable within-image position encoding vector is also added.
- Conveys information about a patch's global position within the image

#### Image sectioning

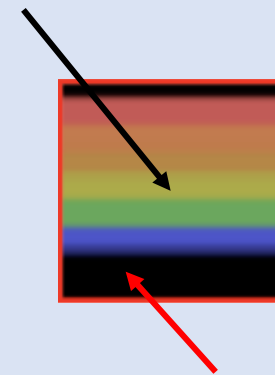


#### Finding coordinate ranges



#### Selecting encoding point

During training, the center point is used:  $x = 64, y = 48$



During evaluation, a random point in the range is used