## Inference accelerator
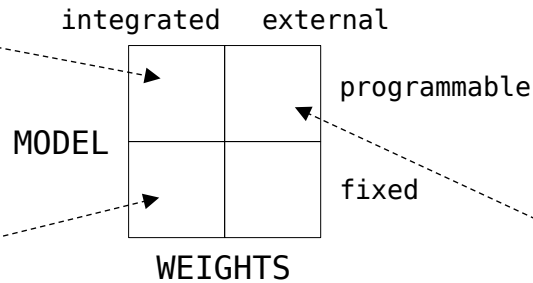
Cerebras
Groq

## Embedded AI

integrated SRAM weights
fixed model
no software stack, weights only

integrated    external

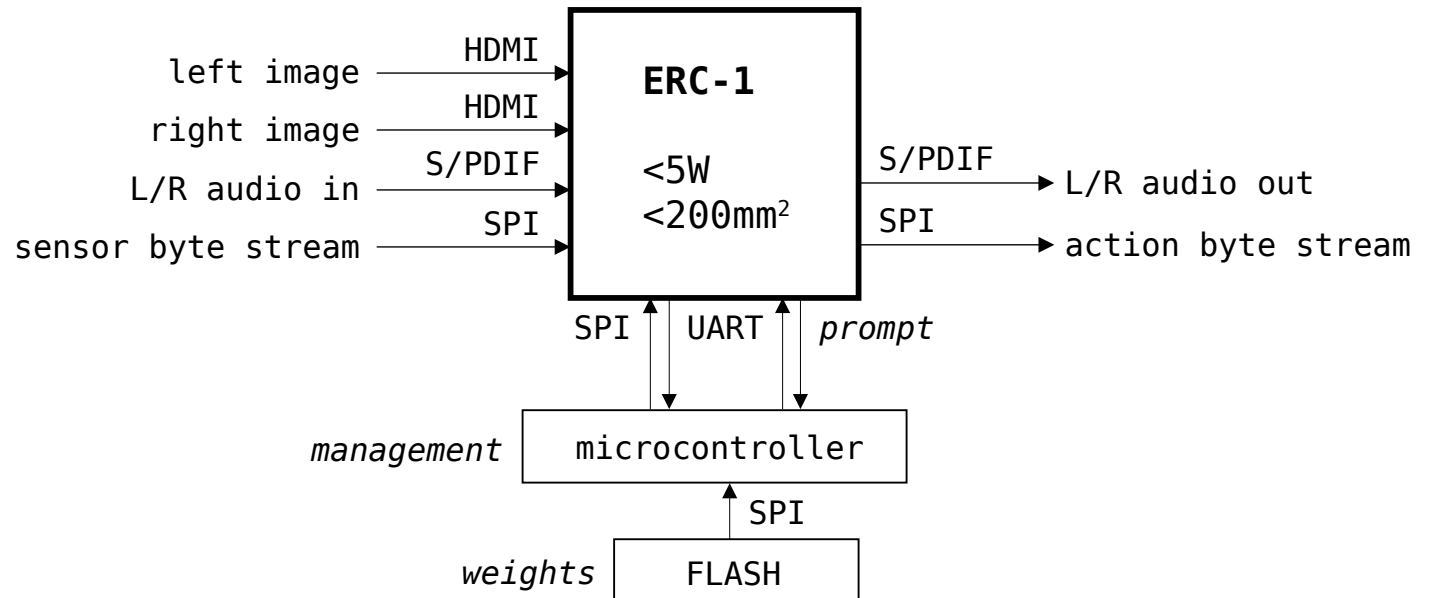MODEL

programmable

fixed

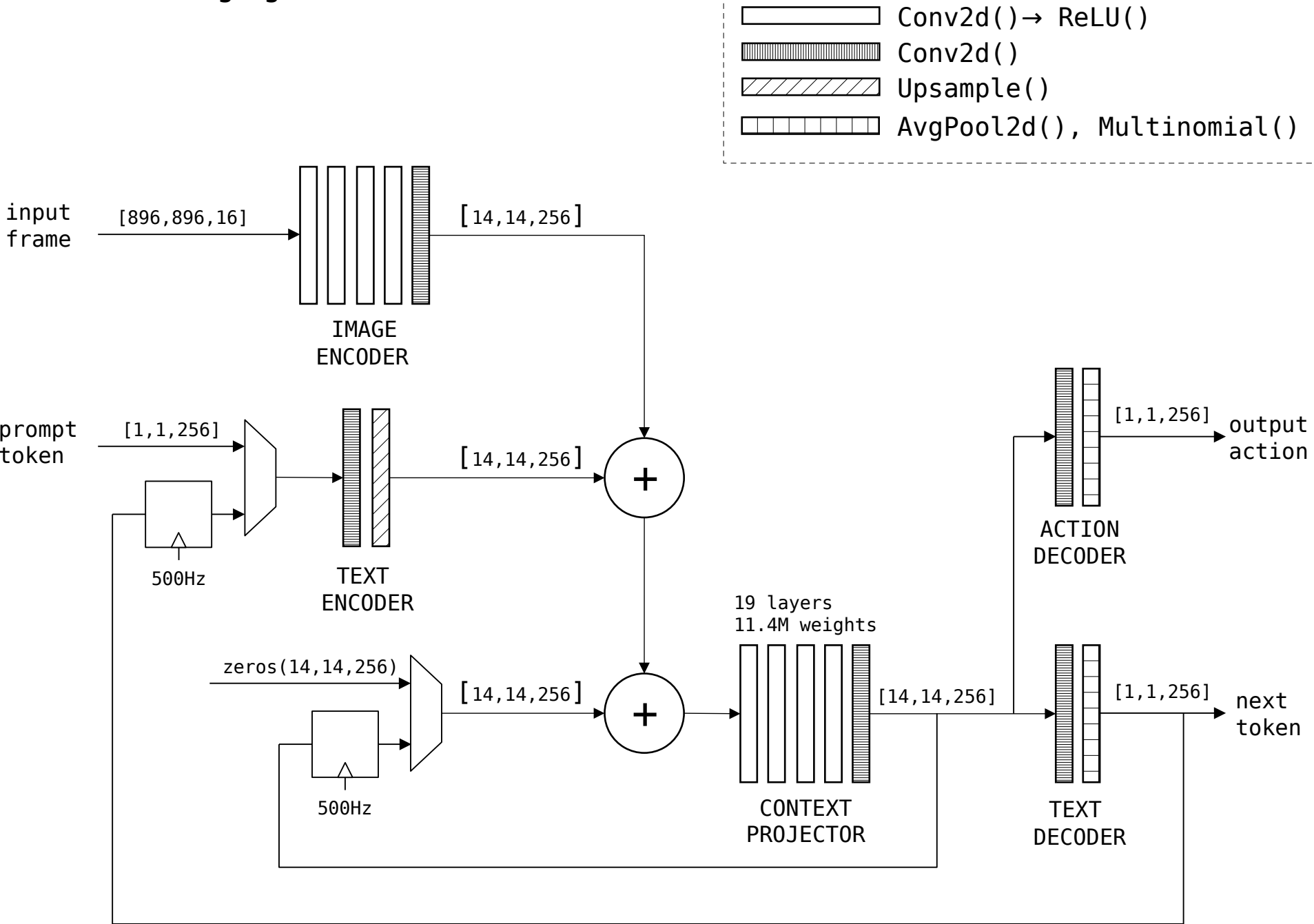WEIGHTS

## AI accelerator (GPU)

external DRAM weights
programmable model
Pytorch software stack

## Embedded Robot Controller

Autonomous robots
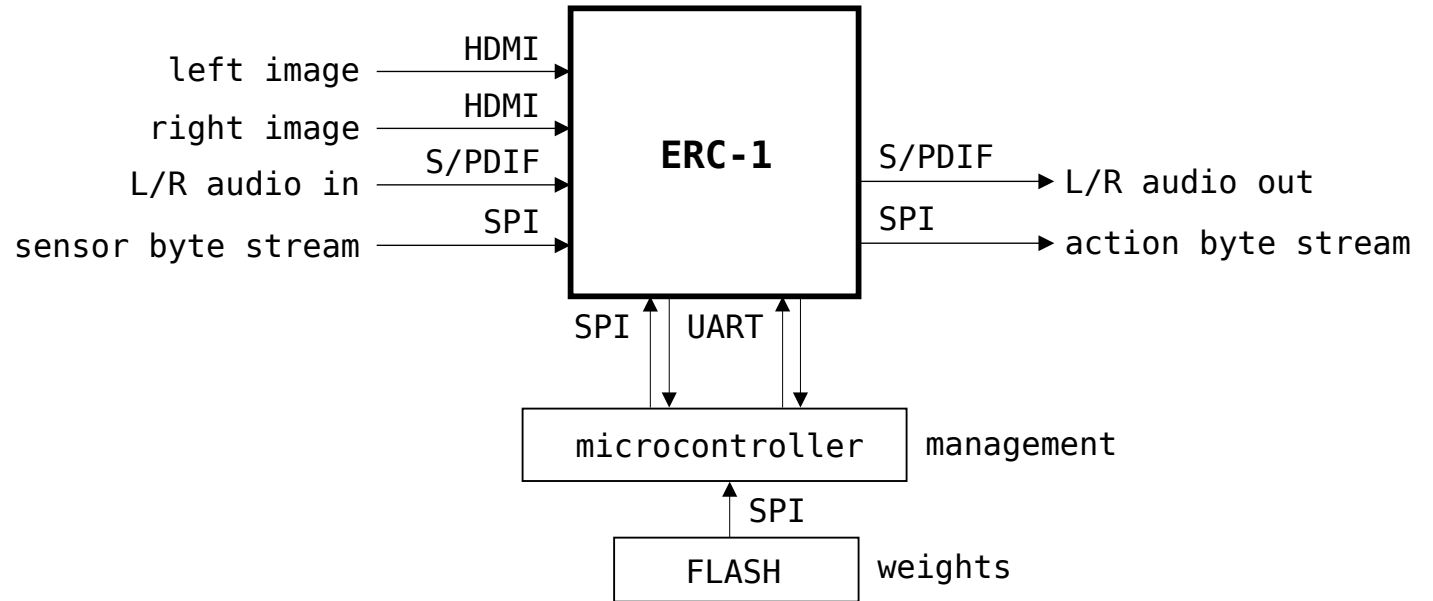Interactive signage
Nuclear fusion reactors

left image — HDMI →

right image — HDMI →

L/R audio in — S/PDIF →

sensor byte stream — SPI →

**ERC-1**

<5W
<200mm$^2$

S/PDIF → L/R audio out

SPI → action byte stream

SPI    UART    *prompt*

*management*    microcontroller

SPI

*weights*    FLASH

# CNN Vision-Language-Action Model

Conv2d()→ ReLU()
Conv2d()
Upsample()
AvgPool2d(), Multinomial()

input frame
[896,896,16]
IMAGE ENCODER
[14,14,256]

prompt token
[1,1,256]
500Hz
TEXT ENCODER
[14,14,256]

+

zeros(14,14,256)
500Hz
[14,14,256]

+

19 layers
11.4M weights
CONTEXT PROJECTOR
[14,14,256]

ACTION DECODER
[1,1,256]
output action

TEXT DECODER
[1,1,256]
next token

# BACKUP

# Embedded Robot Controller

Autonomous robots
Interactive signage
Nuclear fusion reactors

```
left image  ──HDMI──▶ ┌─────────┐
right image ──HDMI──▶ │         │ ──S/PDIF──▶ L/R audio out
L/R audio in ─S/PDIF─▶ │  ERC-1  │
sensor byte stream ─SPI─▶│        │ ──SPI──▶ action byte stream
                       └─────────┘
                        SPI   UART
```

ERC-1

HDMI — left image
HDMI — right image
S/PDIF — L/R audio in
SPI — sensor byte stream

S/PDIF — L/R audio out
SPI — action byte stream

SPI   UART

microcontroller | management

SPI

FLASH | weights

---

**hardened components**
integrated SRAM/ROM weights      ⟺
fixed model

**AI accelerator**
external DRAM weights
programmable model

WEIGHTS

MODEL

programmable

fixed

integrated    external

# CNN Language Model



prompt token

ENCODER

V [E] [H,W,E]

zeros(H,W,C)

[H,W,C]   CAT   [H,W,C+E]

PROJECTOR

[H,W,C]

DECODER

[E]   V   next token

Conv2d()→ ReLU()

Conv2d()

Upsample(), AvgPool2d()

FROZEN

| | |
|---|---|
| V | 50257 |
| E | 256 |
| C | 512 |
| H,W | 8,8 |

# CNN Language Model



prompt token

ENCODER

V [E] [H,W,E]

zeros(H,W,C)

[H,W,C] CAT [H,W,C+E]

PROJECTOR

[H,W,C]

DECODER

[E] V

next token

Conv2d()→ ReLU()

Conv2d()

Upsample(), AvgPool2d()

FROZEN

| V | 50257 |
| E | 256 |
| C | 512 |
| H,W | 8,8 |

# CNN Language Model



prompt token

zeros(H,W,C)

ENCODER

V  [E]  [H,W,E]

[H,W,C]  CAT  [H,W,C]  [E]  V  next token

PROJECTOR  DECODER

Conv2d()→ ReLU()

Conv2d()

Upsample()

FROZEN

| E | 256 |
| V | 50257 |
| C | 256 |
| H,W | 14,14 |

# CNN Language Model



prompt token

zeros(H,W,C+E)

ENCODER

PROJECTOR

DECODER

next token

V   [E]   [H,W,E]

[H,W,C+E]   [H,W,C]   CAT   [H,W,C+E]   [E]   V

Conv2d()→ ReLU()

Conv2d()

Upsample()

FROZEN

| | |
|---|---|
| E | 256 |
| V | 50257 |
| C | 256 |
| H,W | 14,14 |

# CNN Vision-Language Model

*caption*

prompt

V  [256]  [H,W,256]

TEXT ENCODER

image  [896,896,3]
spectrogram  [896,896,1]

[H,W,128]

IMAGE ENCODER

C

zeros(H,W,640)

[H,W,640]  [H,W,256]

CONTEXT MAP

C  [H,W,640]

[256]  V

TEXT DECODER

Conv2d()→ ReLU()
Conv2d()
Upsample()→ Conv2d()

| H,W | 32,32 |

# CNN Perception-Language-Action Model

*inner caption*

prompt

V [256]

[H,W,256]

TEXT ENCODER

left eye [896,896,3] [H,W,128]

right eye [896,896,3] [H,W,128]

spectrogram [896,896,3] [H,W,64]

IMAGE ENCODERS

[256] V

TEXT DECODER

zeros(H,W,C)

[H,W,832]

[H,W,256]

[H,W,832]

CONTEXT MAP

C

pose
map

POSE DECODER

Conv2d()→ BatchNorm2d()→ ReLU()

Conv2d()

Upsample()→ Conv2d()

H,W    32,32

**PHASE 1:**
CNN Perception-Language-Action model (CNN-PLA), running on GPU
>20fps, <100ms latency, <400W

**PHASE 2:**
CNN-PLA model running on FPGA
>100fps, <15ms latency, <100W

**PHASE 3:**
CNN-PLA model running on ASIC at >300fps, <5ms latency, <10W

Seated humanoid :
1-DoF spine
2-DoF shoulder
2-DoF neck
6-DoF arms
6-DoF hands

Data collection :
VR teleoperated
spine/shoulder : swivel chair with up/down control
neck : VR gyroscope
arms/hands : TBD

ROBOT
896x896 RGB at 120Hz, left and right eye
896x896 spectrogram at 120Hz, audio and tactile signals
29-DoF target pose at 120Hz

TELEOPERATOR
transcribed verbal commentary at 120 characters/s
synthetically captioned using open source VL model, 1 character/frame

# CNN Perception-Language-Action Model

*inner caption*

prompt

V [256]

[H,W,256]

**TEXT ENCODER**

left eye [896,896,3] [H,W,128]

right eye [896,896,3] [H,W,128]

spectrogram [896,896,3] [H,W,64]

**IMAGE ENCODERS**

C

C

C

[256] V

**TEXT DECODER**

zeros(H,W,C)

[H,W,832] [H,W,256] C [H,W,832]

**CONTEXT MAP**

pose map

**POSE DECODER**

Conv2d()→ BatchNorm2d()→ ReLU()

Conv2d()

Upsample()→ Conv2d()

H,W    32,32

**CNN Language Model**

prompt token

zeros(H,W,C)

V  [E]

ENCODER CNN

[H,W,C]

PROJECTOR CNN

[H,W,C]

[H,W,C]

+

[H,W,C]

DECODER CNN

[E]  V

next token

Conv2d()→ ReLU()

Conv2d()

Upsample()

AvgPool()

| E | 256 |
| V | 50257 |
| C | 384 |
| H,W | 14,14 |

Devices

- TE256A      : Token Encoder, vocab_size=256, ASCII
- TD256A      : Token Decoder, vocab_size=256, ASCII
- CP384V      : Context Projector, n_embd=384, VGG
- IE384R      : Image Encoder, n_embd=384, Resnet
- TD
-

alt = lite-base-resnet-xeno

*inner voice*

prompt

V  [C]  [H,W,C]

TEXT ENCODER

perception  [768,768,3]  [H,W,C]

IMAGE ENCODER

+

TEXT DECODER

[C]  V

zeros(H,W,C)  [H,W,C]  [H,W,C]

CONTEXT PROJECTOR

+

POSE DECODER

[C]  pose

Conv2d(kernel=3)→ ReLU()

Conv2d(kernel=1)

Upsample(scale=3)→ Conv2d(kernel=3)→ ReLU()

| V | 256 |
| C | 384 |
| H,W | 81,81 |

alt2



prompt token

V [C] [H,W,C]

REPLICATE

zeros(H,W,C)

[H,W,C] [H,W,C] [H,W,C] [C] V

next token

+

CONTEXT CNN

DECODER CNN

Conv2d(kernel=3)→ BatchNorm2d()→ ReLU()

Conv2d(kernel=1)

Upsample()

| | |
|---|---|
| V | 256 |
| C | 384 |
| H,W | 81,81 |

alt1

prompt
token

V   [C]   [H,W,C]

TEXT ENCODER CNN

zeros([H,W,C])

[H,W,C]   [H,W,C]   +   [H,W,C]   [C]

STATE TRAJECTORY CNN

TEXT DECODER CNN

alt1

current token $\xrightarrow{\text{V}}$ [C] → TEXT ENCODER CNN → [H,W,C]

zeros([H,W,C]) → [H,W,C] → CONTEXT TRAJECTORY CNN → [H,W,C] → + → [H,W,C] → TEXT DECODER CNN → [C] →

alt1

current token —V→ [CNN ENCODER] —[C]→ [H,W,C]

CNN ENCODER

current context —[H,W,C]→ [CNN PROJECTOR] —[H,W,C]→

CNN PROJECTOR

(+) —[H,W,C]→ [CNN DECODER] —[C]→ V→

CNN DECODER

[H,W,C]

multimodal alt1

current token —V→ ‖ —[C]→ ▯ CNN ENCODER —[H,W,C]→ (+) —[H,W,C]→ CNN DECODER —[C]→ ‖ —V

CNN ENCODER

CNN DECODER

[H,W

current frame —[768,768,8]→ CNN ENCODER —[H,W,C]→ (+)

current context —[H,W,C]→ CNN PROJECTOR —[H,W,C]→

multimodal alt1

current frame $\xrightarrow{\text{[768,768,8]}}$ CNN ENCODER $\xrightarrow{\text{[H,W,C]}}$

CNN ENCODER

current token $\xrightarrow{V}$ [C] $\rightarrow$ CNN ENCODER $\xrightarrow{\text{[H,W,C]}}$

CNN ENCODER

$+$ $\xrightarrow{\text{[H,W,C]}}$ CNN DECODER $\xrightarrow{\text{[C]}}$ V

CNN DECODER

[H,W

np.zeros([H,W,C]) $\rightarrow$ $\xrightarrow{\text{[H,W,C]}}$ CNN ENCODER $\xrightarrow{\text{[H,W,C]}}$

CNN ENCODER

multimodal alt1

current frame [768,768,16] → CNN ENCODER → [H,W,C]

[H,W,C] → CNN ENCODER → [768,768,1...

current token V → [C] → CNN ENCODER → [H,W,C]

+

[H,W,C] → CNN DECODER → [C] → V

np.zeros([H,W,C]) → [H,W,C] → CNN ENCODER → [H,W,C]

[H,W

alt1



current token —int→ [ ] —C→ CNN ENCODER —H,W,C→ (+) —H,W,C→ CNN DECODER —C→ [ ] i

current context —H,W,C→ CNN PROJECTOR —H,W,C→

H,W

current token — int — C — CNN ENCODER — H,W,C — + — H,W,C — CNN DECODER — C — int — ne to

current context — H,W,C — × — ne co

H,W,C

γ (gamma)

alt2

current token →[int]→ [C] → CNN ENCODER → H,W,C

current context →[H,W,C]→ CNN PROJECTOR → H,W,C

CAT → H,W,2*C → H,W,C → CNN DECODER → C

H,W

alt3

int

current
token

E

CNN ENCODER

H,W,C

CAT

H,W,2*C

H,W,
C

CNN PROJECTOR

H,W,C

CNN DECODE

H,W,C

current
context

current
token →int→ C →| CNN ENCODER →H,W,C→ (+) →H,W,C→ CNN PROJECTOR →H,W,C→ CNN DECODER → C →

current
context H,W,C → (×) 

γ
gamma

current
token      int        C      H,W,C        H,W,C        C      int     ne
                                                               to

CNN ENCODER        +       CNN DECODER

current    H,W,C                                              H,W,C   ne
context                                                               co

current token → int → C → CNN ENCODER → H,W,C → (+) → H,W,C → CNN PROJECTOR → H,W,C → CNN DECODER → C →

current context → H,W,C

current token →[int]→ C →□ CNN ENCODER → H,W,C

CNN ENCODER

current context → H,W,C → CNN PROJECTOR → H,W,C → (+) → H,W,C → CNN DECODER → C →[ i

CNN PROJECTOR

CNN DECODER

H,W