→ segmentation has always been done using a FCN leveraging an encoder-decoder architecture, This approach was very successful, until it was clear that The bottleneck of the Model was context in very precise segmentation Tasks ( self driving, inventory... )

⟹ researchers Tried To augment context by increasing receptive field, but receptive field increased linearly with [number of conv layers × kernel size]

● increasing number of conv layers ⟹ More depth ⟹ More compute + risk of overfitting

● increasing kernel size ⟹ Model won't be able To detect intricate details which is especially important in segmentation Tasks, due To pixel-level labeling //

⇒ receptive field remains unsolved until Now

→ SETR comes To solve This receptive field problem with Transformer architecture which enables context To The furthest extent known Till now ⫽

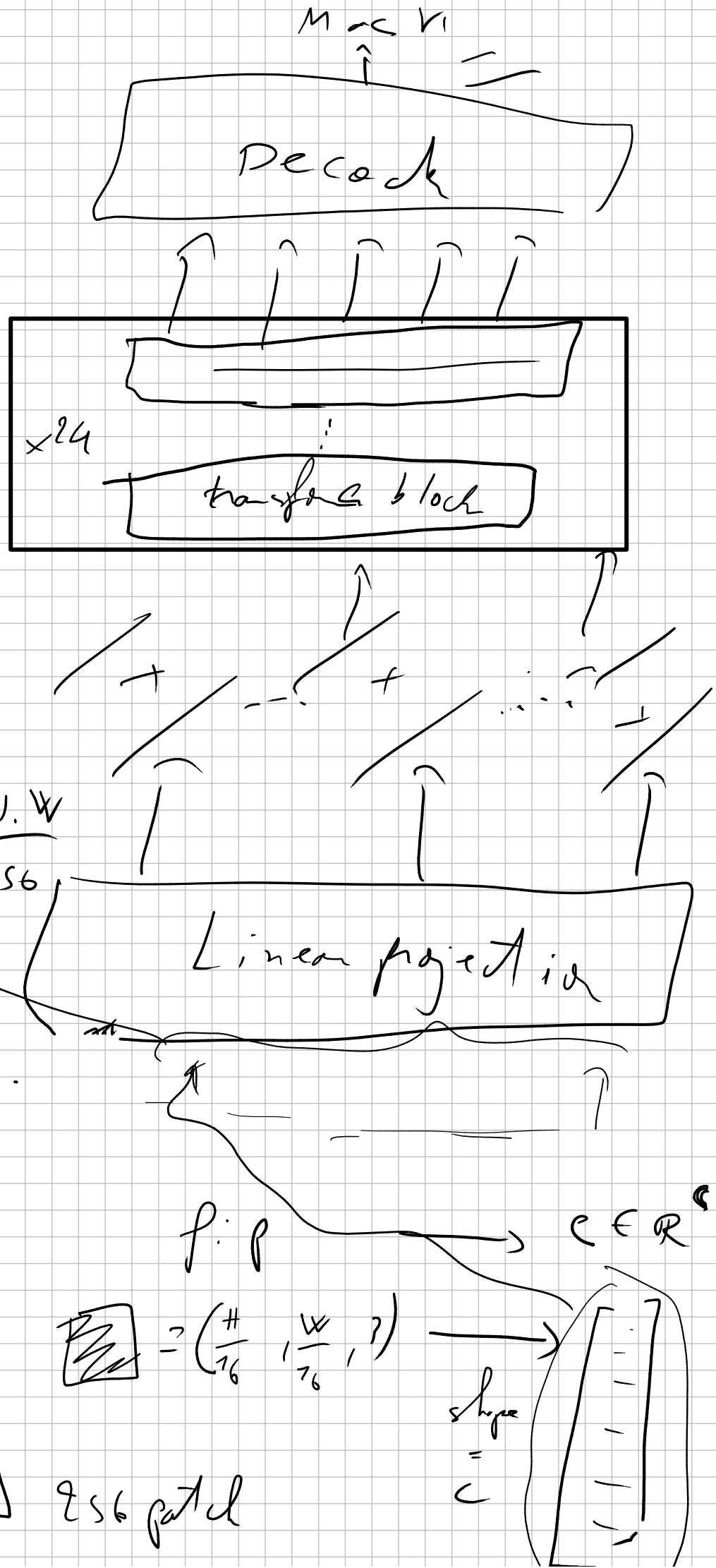3 main architecture proposed in The paper:

SETR
MLA

SETR
segmentation

SETR
PUP

Transformer

multi-level feature aggregation

progressive up sampling

# SETR:

Max $v_1$

Decode

$\times 24$

transformer block

$+$ $+$ $\cdots$

sequence length $N = L = \dfrac{H \cdot W}{256}$

Linear projection

$f : P$

$\in \mathbb{R}^c$

$16 \{$

$\boxed{\phantom{x}} \to \left( \dfrac{H}{16}, \dfrac{W}{16}, 3 \right) \longrightarrow$

shape $=$ $c$

256 patch

## decoder design:

$$\text{decoder input} = \left[ \overbrace{\left[\phantom{|||}\right], \left[\phantom{|||}\right], ---, }^{16} \right]$$

$\in \mathbb{R}^{E}$

$L = A$

$\dfrac{H \cdot W}{256}$ data

$$\text{decoder output} = \left[ \cancel{\varnothing}_j, \cancel{\varnothing}, \qquad\qquad \right]$$

$\in \mathbb{R}^{\frac{H}{16} \times \frac{W}{16} \times C}$

## (1) Naïve upsampling:

$\times 2$

$X \in \mathbb{R}^{C} \begin{array}{c} C \\ C \end{array} \begin{cases} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{cases} = \left[ \begin{array}{|c|c|c|c|} \hline & 1 & & \\ \hline \end{array} \right]$

$1 \times 1 + \text{batch norm}$
conv
$19$ [nb of classes]
filters

$$1 \times 1 \text{ cov} + \text{batch Nor} + 1 \times 1 \text{ cov}$$

(right brace) num classes

## (2) PUP:

$$\frac{H \cdot W}{256} = \text{number of embedding} \quad \overset{as}{\text{input to decode}}$$

$$C = \text{dimension of embedding}$$



Reshape

$$\frac{H}{16} \times \frac{W}{16} \times 1024$$

$$\frac{H \cdot W}{256} \times C \quad \text{us hyperd}$$

scenario

patch was reshaped at the start

conv
→ 2x

transpose conv (kernel size(?, ?))

$$\frac{H}{8} \times \frac{W}{8} \times 256 \longrightarrow \quad conv \atop 2x \quad \frac{H}{4} \times \frac{W}{4} \times 256$$

$$H \times W \times 19 \quad \longleftarrow \quad conv\, 2x$$

= Mask

→ We do this since one-step upscaling introduced in naive decoder might introduce Noise as it might underfit The scaling process

(!) SETR-PUP have more params Then SETR-Naive because of this deeper decoder
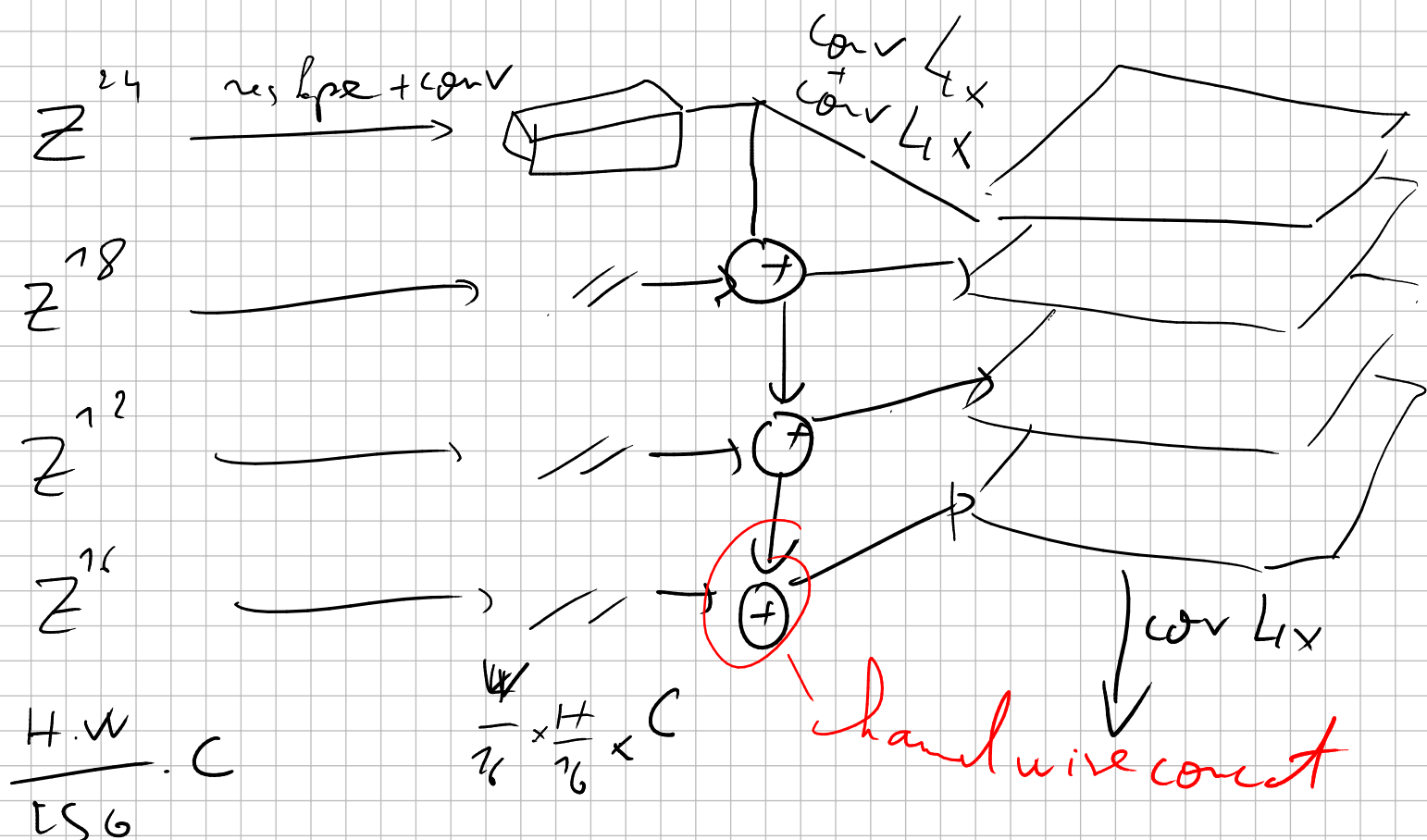
# (3.) MLA :

→ This decoder is even longer ( Params ++)

further proving scale MATTERS , as This

decoder out performed (2) & (1)

⟹ Multi-level feature aggregation :

$$Z^m = \text{output of } n^{\underline{R}} \text{ Transformer blocks}$$

$$\text{in The encoder} : \left( \frac{H \cdot W}{256} , C \right)$$

$Z^{24}$  $\xrightarrow{\text{reshape +conv}}$

$Z^{18}$

$Z^{12}$

$Z^{6}$

conv 4x + conv 4x

conv 4x

conv 4x

channel wise concat

$\frac{H \cdot W}{256} \cdot C$          $\frac{W}{16} \times \frac{H}{16} \times C$

⟹ we fuse all feature from The Transformer layers

$\Rightarrow$ SETR-MLA

(!) $\text{Paras}_{MLA} > \text{Paras}_{PUP} > \text{Paras}_{Naive}$

$\leftarrow$ Scale

◯ in the paper ViT & DEiT
weights were used + embeddings
to make up for data - shortage

$\Rightarrow$ ViT : image understanding
( long - range context )

finetuned for segmentation

Task using 3 different

decoders

$\Rightarrow$ achieves $\begin{cases} \cdot \text{ SOTA on Pascal} \\ \cdot \text{ Competitive results on cityscapes} \\ \cdot \text{ 1}^{st} \text{ on ADE20k} \end{cases}$