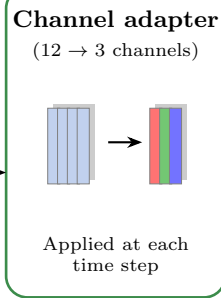
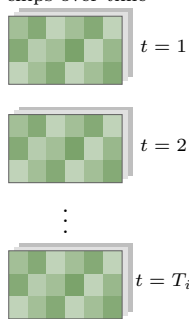


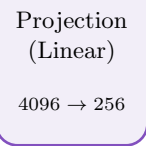
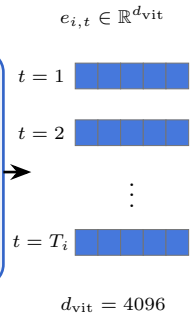
Image sequence

$$X_i = \{x_1, \dots, x_{T_i}\}$$

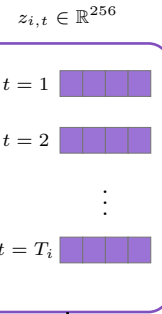
12-band Sentinel-2
chips over time



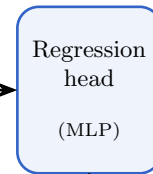
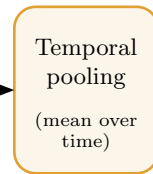
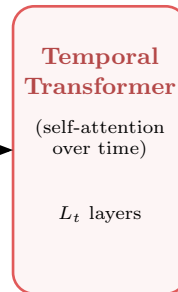
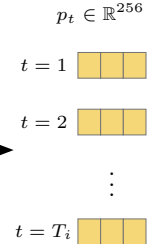
Per-time-step image embedding



Projected per-time-step image embeddings



Learned temporal positional embeddings



Scalar yield
prediction

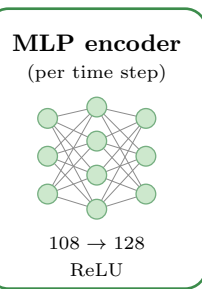
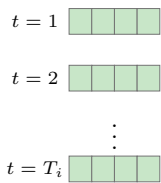
$$\hat{y}_i$$

(t/ha)

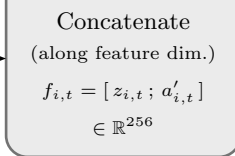
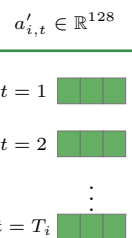
Auxiliary feature sequence

$$A_i = \{a_1, \dots, a_{T_i}\}$$

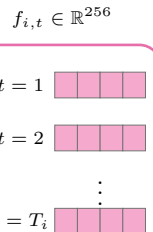
$$a_{i,t} \in \mathbb{R}^{108}$$



Per-time-step auxiliary embeddings



Fused per-time-step features



- T_i : number of time steps (field-specific)
- H, W : spatial height and width of chips
- N_p : number of patches per image
- d_{vit} : ViT embedding dimension (4096)
- 128 : projection / embedding dimension
- L_t : number of temporal transformer layers