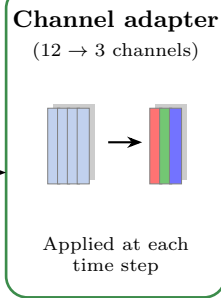
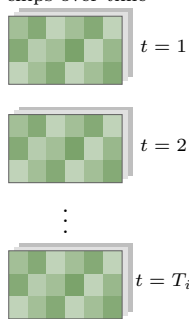


Image sequence

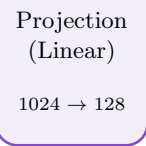
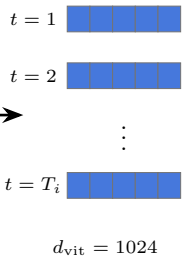
$$X_i = \{x_1, \dots, x_{T_i}\}$$

12-band Sentinel-2
chips over time



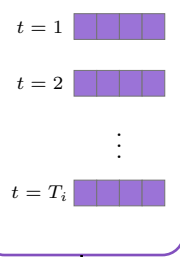
Per-time-step image embedding

$$e_{i,t} \in \mathbb{R}^{d_{vit}}$$



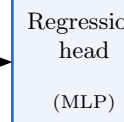
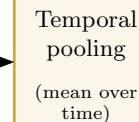
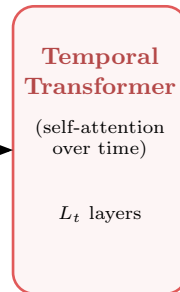
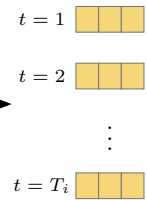
Projected per-time-step image embeddings

$$z_{i,t} \in \mathbb{R}^{128}$$



Learned temporal positional embeddings

$$p_t \in \mathbb{R}^{128}$$



Scalar yield
prediction

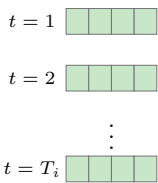
$$\hat{y}_i$$

(t/ha)

Auxiliary feature sequence

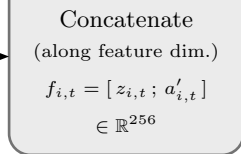
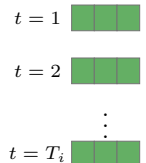
$$A_i = \{a_1, \dots, a_{T_i}\}$$

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



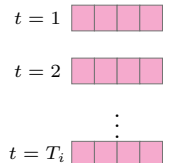
Per-time-step auxiliary embeddings

$$a'_{i,t} \in \mathbb{R}^{128}$$



Fused per-time-step features

$$f_{i,t} \in \mathbb{R}^{256}$$



- 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 (1024)
- 128 : projection / embedding dimension
- L_t : number of temporal transformer layers