# Quiz

**Quiz 1:** Match each feature with the correct model (CNN or ViT).

| Feature | Model |
|---|---|
| (A) Uses Self-attention mechanism (Transformers) to process images | _ViT_ |
| (B) Use convolutional filter slides across the image and extracts features | _CNN_ |
| (C) Treat images as sequences of patches | _ViT_ |
| (D) Works well even with small datasets | _CNN_ |
| (E) Needs large datasets | _ViT_ |
| (F) Requires positional encodings to retain spatial information | _ViT_ |

**Quiz 2:** In patch embedding process, if the flattened 2D patches is $x_p = \begin{pmatrix} 1 & 2 \\ 1 & 1 \end{pmatrix}$

and the linear projection matrix is $\mathbf{E} = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 4 & 4 \end{pmatrix}$, what is $x_p^2 \mathbf{E}$ after patch embedding?

Hint: $z_0 = [x_{class}; \; x_p^1 \mathbf{E}; \; x_p^2 \mathbf{E}; \ldots; \; x_p^N \mathbf{E}] + \mathbf{E}_{pos}$

$$x_p^2 \mathbf{E} = \begin{pmatrix} 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 \\ 1 & 4 & 4 \end{pmatrix} = \begin{pmatrix} 2 & 5 & 5 \end{pmatrix}$$