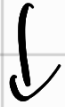


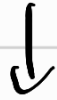
→ Giống như vẽ Bài toán Object Classification



Sự kết hợp giữa CV và NLP

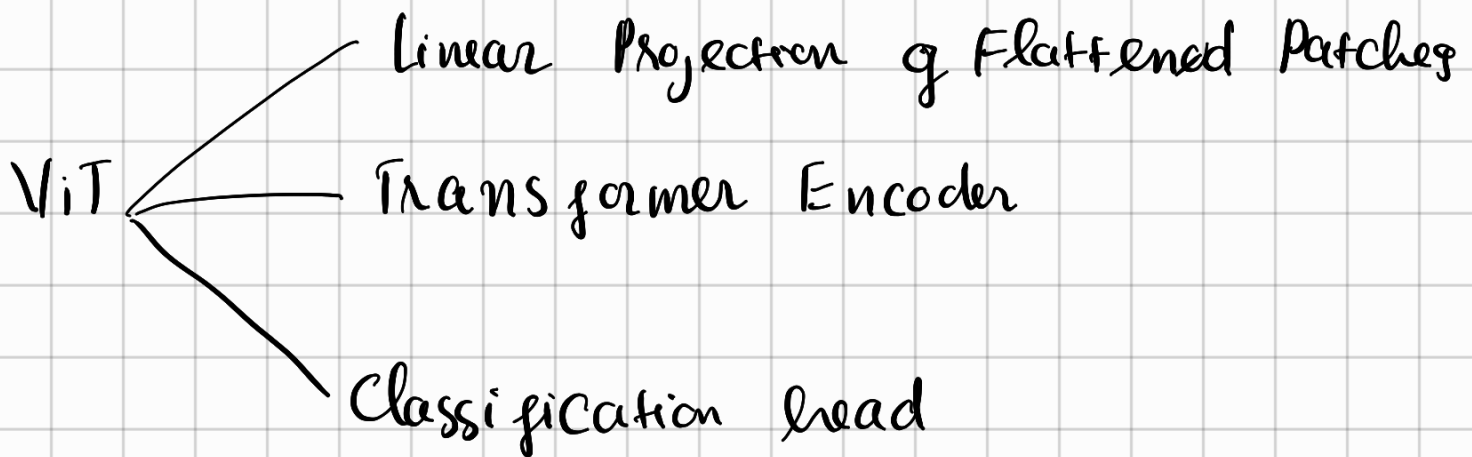


Visual Transformer (ViT)



Là sự kết hợp giữa 1 phần kiến trúc của

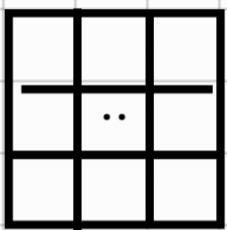
Transformer và các khối MLP
(Multilayer Perceptron)



④ Linear Projection of Flattened Patches

2.1 Patch Embedding

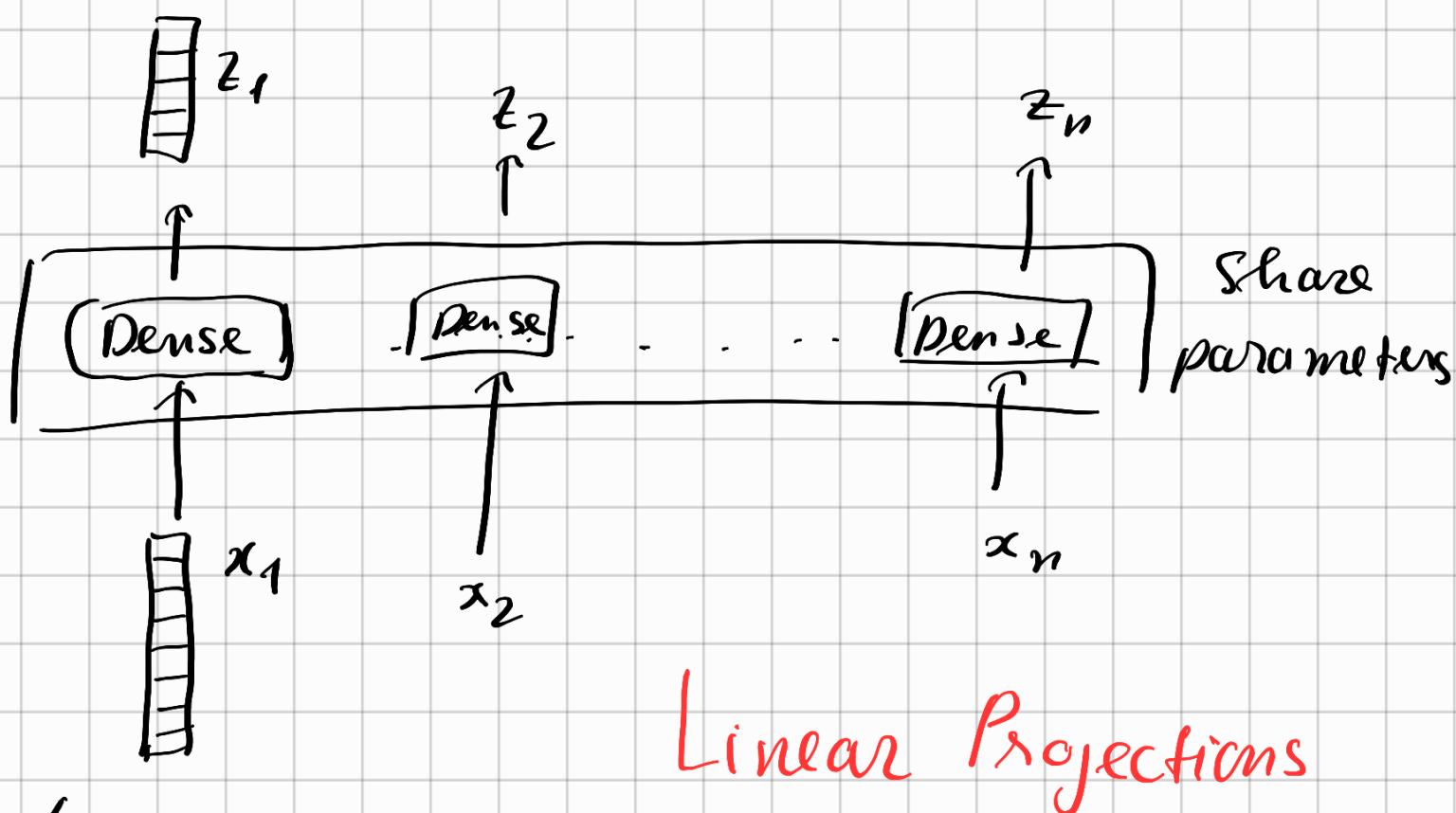
↳ chia ảnh ra làm các patch nhỏ



split → 9 patch.

image

→ Dưa cái patch về dạng vector = flattend
các patches ra
- mô hình



Linear Projections

→ bản chất : là một lớp Dense vs input là flattend vector của các batch

→ Output: embedding vector
tương ứng vs từng batch

$$z_i = W * x_i + b \quad \left| \begin{array}{l} x_i: \text{flattened vector} \\ z_i: \text{output of } x_i \\ W: \text{matrix embedding} \end{array} \right.$$

2.27 Positional Embedding

→ thêm thông tin về vị trí cho mỗi patch ở đầu z_i

2.31 Class Embedding

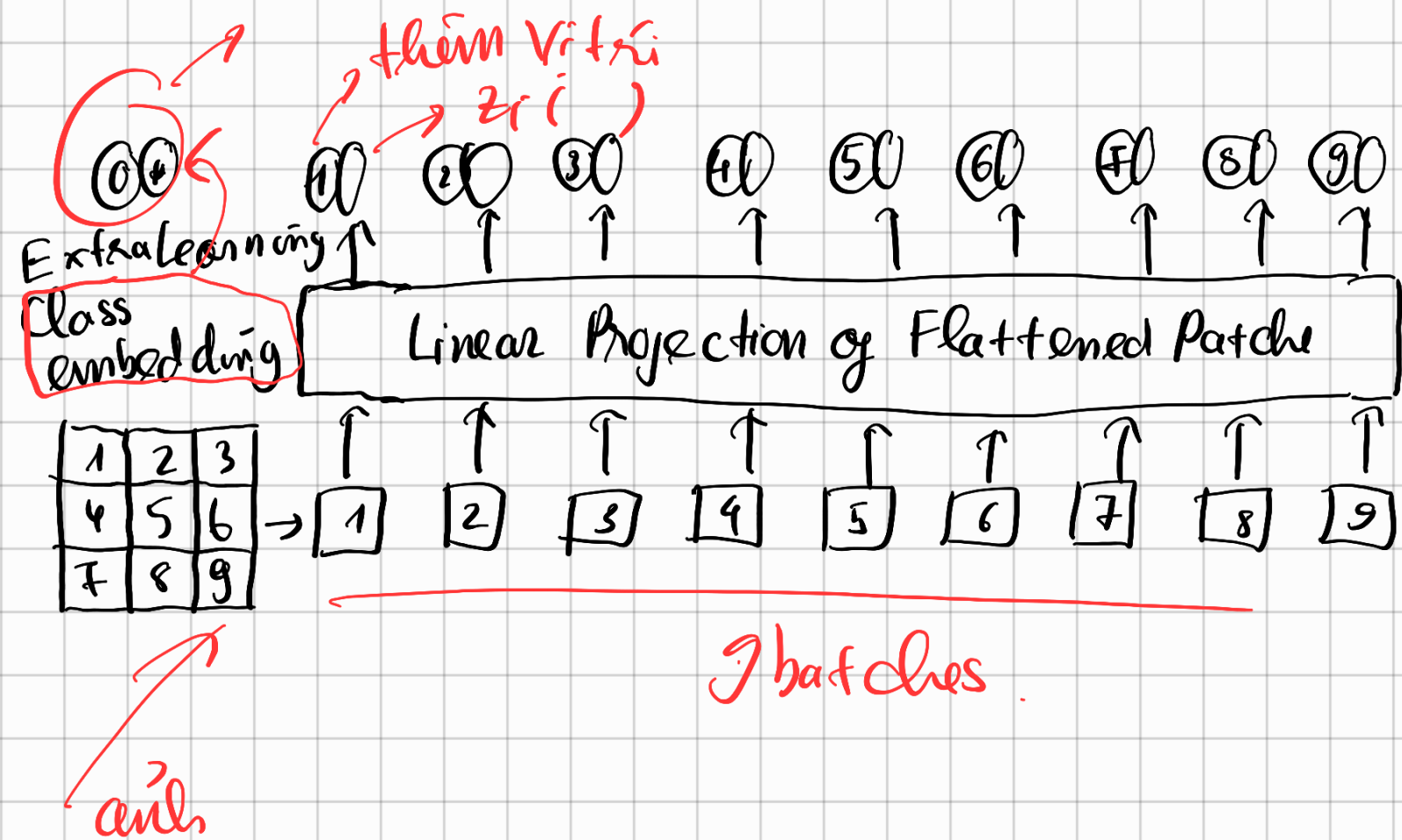
→ ý tđ là mã hoá class

→ giúp mô hình "hiểu" hơn về các output class label.

⊕ Tổng quan :

Sau 8 step này từ một bức ảnh gốc sử dụng ảnh được chia là 3×3

→ mô hình ở dưới nhé



2.2 Transformer Encoder
 + 7 Self Attention layer
 + 7 Multi head Attention

2.3 Classification Head.

một khối MLP (Multilayer Perceptron)

Input
context vector c

Suất ra tổng cộng
vs các class

31 Training Strategies

3 bước training

+? Khởi tạo Pre-training

+? Fine-tune

+? Testing

