

## Tiền huấn luyện mô hình

$$\mathcal{D} = \{(I_w, \textcolor{red}{T}_w)\} + \{(I_h, \textcolor{green}{T}_h)\}$$

Tiền huấn luyện

Hỗn hợp các bộ  
mã hóa - giải mã  
đa phương thức

Các tác vụ chuyên biệt

Bộ lọc  
(Bộ mã hóa văn bản  
có điều kiện hình ảnh)

Tinh chỉnh  
ITC&ITM

$$\{(I_h, \textcolor{green}{T}_h)\}$$

Tinh chỉnh  
LM

Bộ chú thích  
(Bộ giải mã văn bản  
có điều kiện hình ảnh)

## Tạo dữ liệu mới

Lọc

$$\{(I_w, \textcolor{red}{T}_w)\}$$

$$\{I_w\}$$

Chú thích

$$\{(I_w, \textcolor{green}{T}_w)\} + \{(I_w, \textcolor{red}{T}_s)\}$$

$$\begin{aligned}\mathcal{D} = & \{(I_w, \textcolor{green}{T}_w)\} \\ & + \{(I_w, \textcolor{red}{T}_s)\} \\ & + \{(I_h, \textcolor{green}{T}_h)\}\end{aligned}$$

luồng nạp dữ liệu vào  
mô hình

luồng xử lý dữ liệu

$I_w$  hình ảnh từ mạng

$I_h$  hình ảnh được người  
gán nhãn

$T_w$  văn bản từ mạng

$T_w$  văn bản từ mạng sau  
khi lọc

$T_s$  văn bản tổng hợp

$T_s$  văn bản tổng hợp sau  
khi lọc

$T_h$  văn bản được người  
gán nhãn

Màu **đỏ** là dữ liệu trước lọc  
nhiều, màu **xanh** là dữ liệu  
sau khi lọc nhiều.