

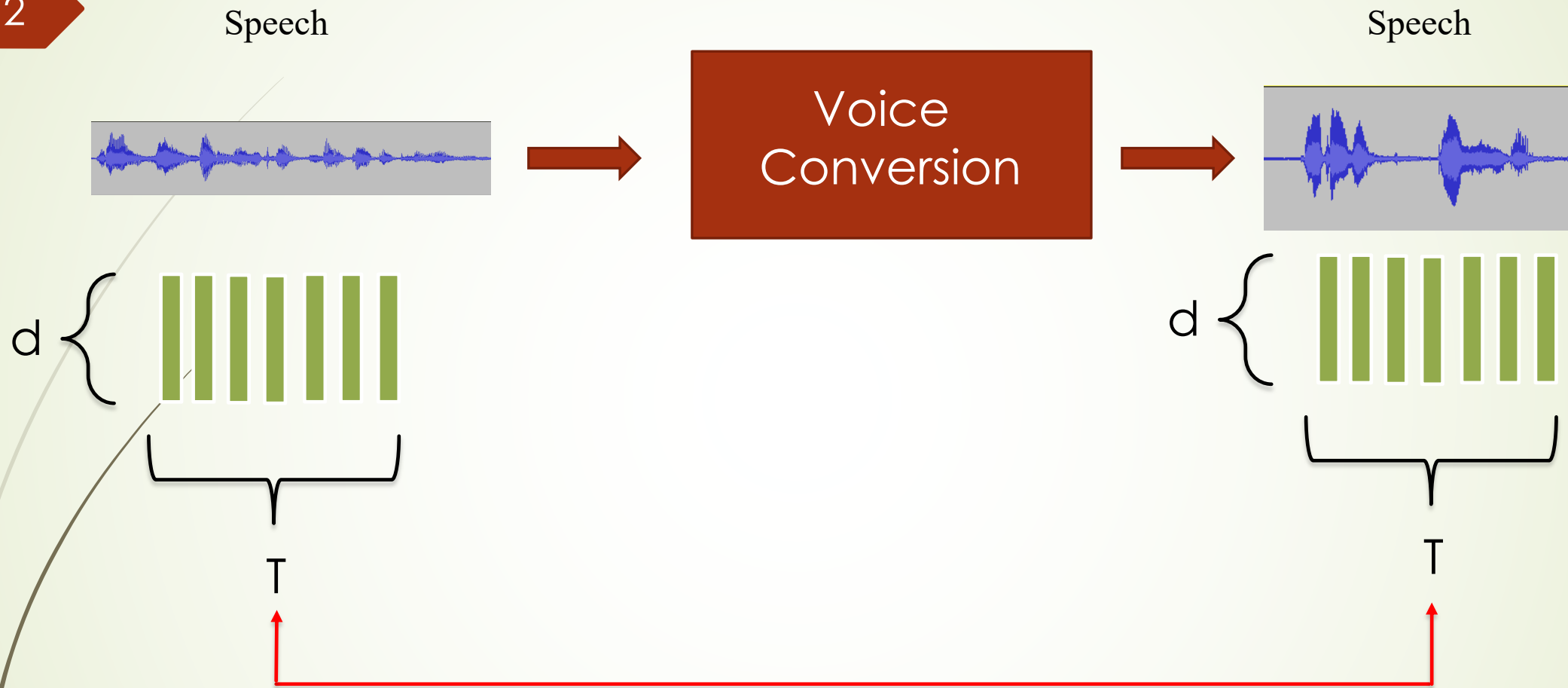
Multi Target voice conversion and cross- language

1

長庚大學 資工所 劉祈宏
指導教授 呂仁園

What is Voice Conversion

2



What can be preserved?

What is changed?

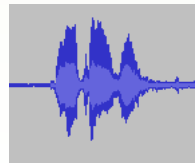
內容

許多都可以，例如語者



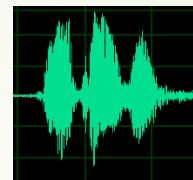
Voice Conversion 技術發展分類

with 平行語料



你好

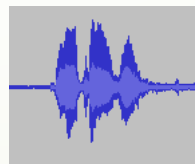
source



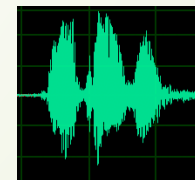
你好

target

without 平行語料



你好

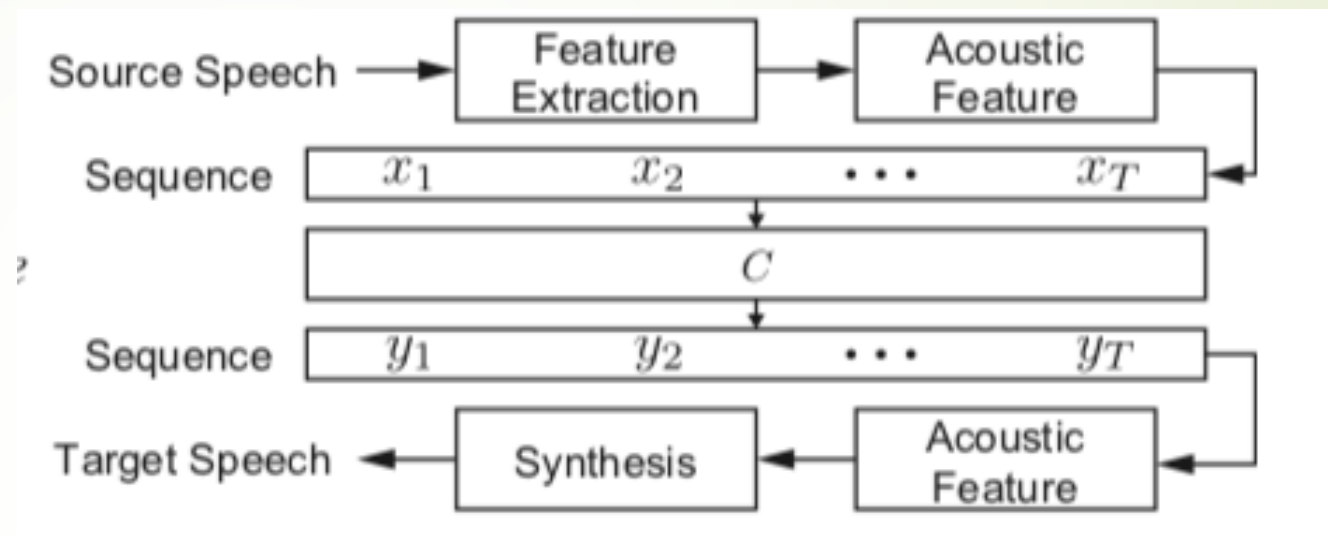


天氣不錯

平行語料

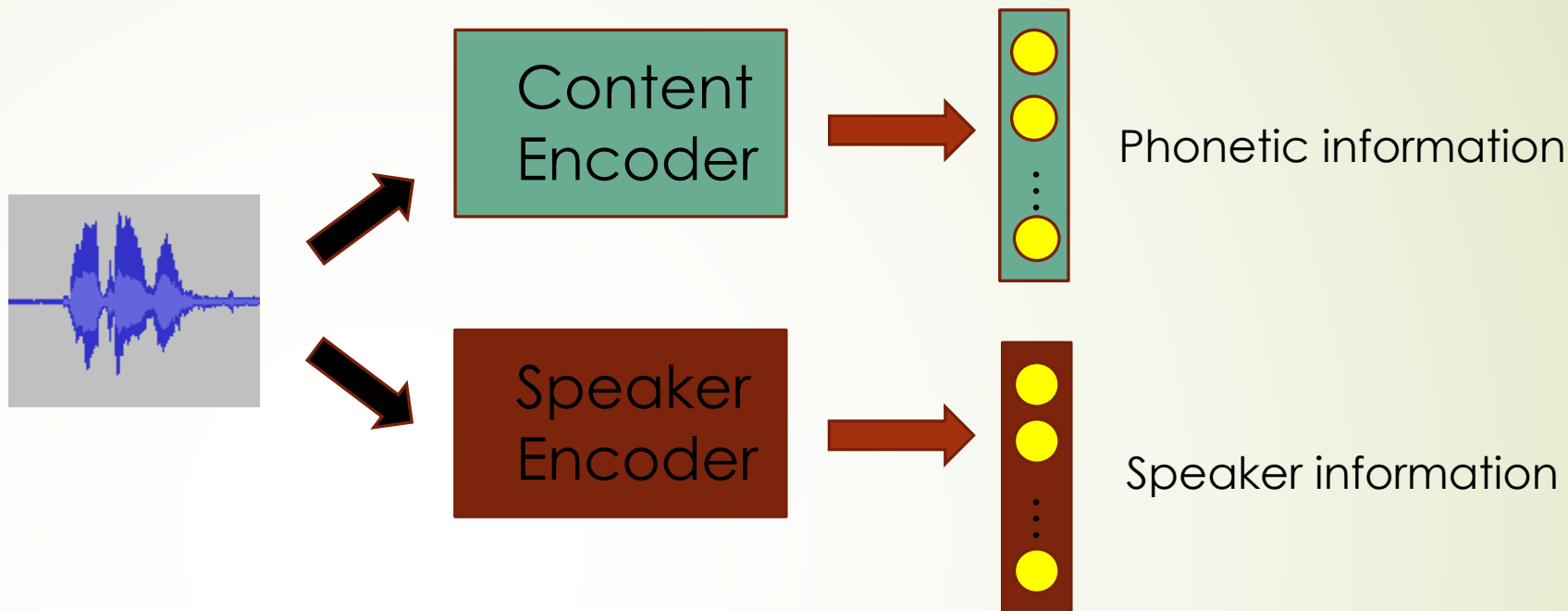
非平行語料

Sequence to sequence model



6

平行語料



非平行語料

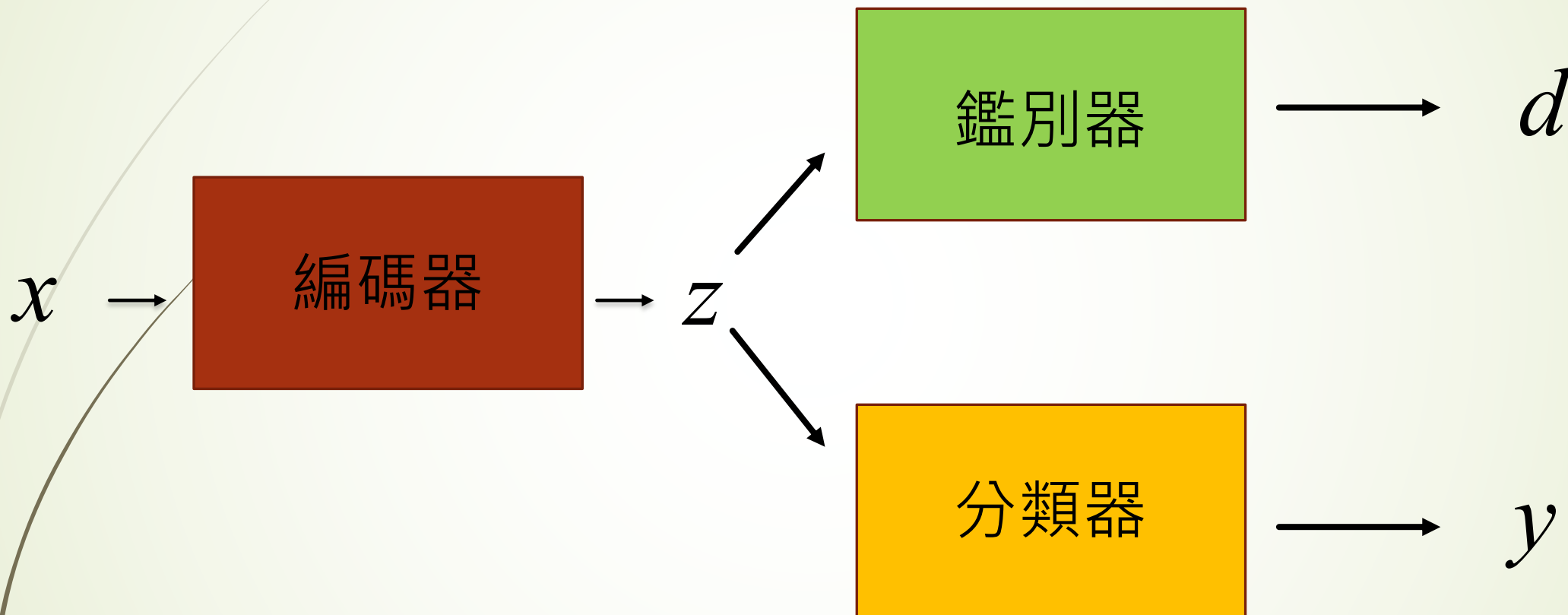
Feature Disentangle

(特徵解纏)

Direct Transformation (直接轉換)

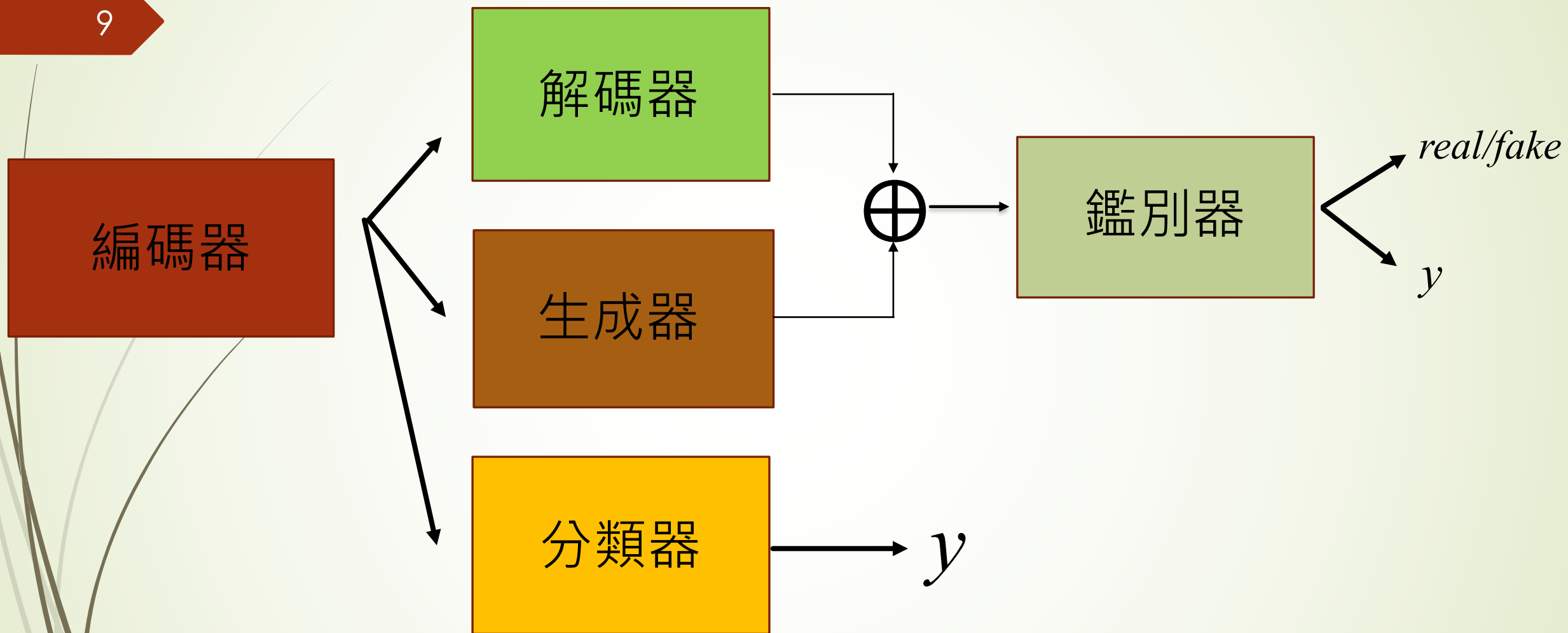
使用對抗訓練方法進行解纏特徵學習

7



x 為資料， y 為分類器的輸出， z 為學習之潛在特徵

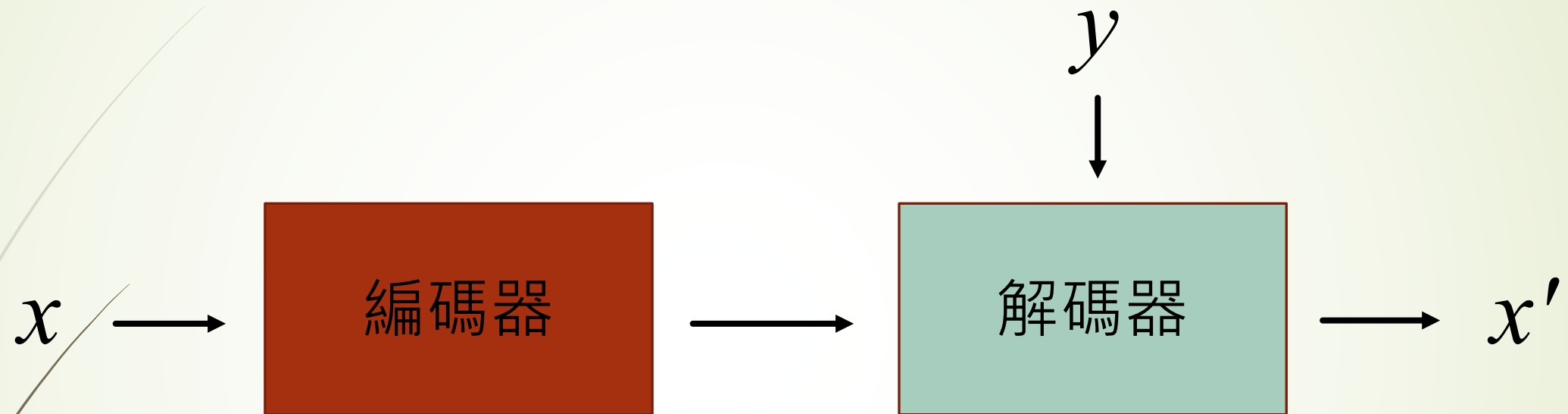
如果只有非平行語料，又想實現多目轉換



其中 y 為語者編號， \oplus 為個別元素相加(Elementwise Addition)。

模型架構與訓練過程

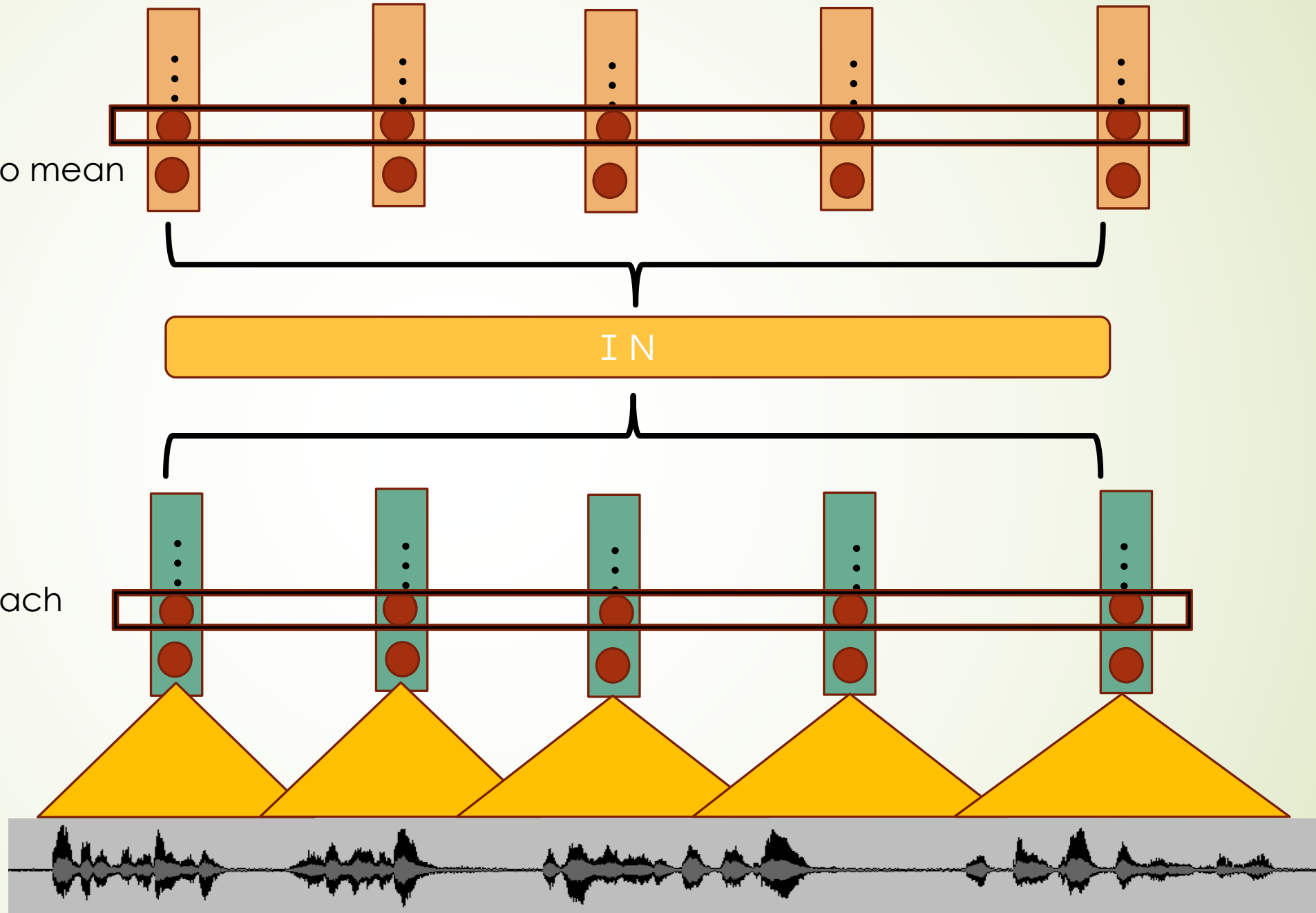
- ➡ 預訓練階段:此階段如同訓練一個自編碼器。
- ➡ 解纏特徵學習階段:此階段會使用一個輔助分類器來幫助編碼器壓縮出不含語者資訊的特徵。
- ➡ 生成對抗網路階段:此階段的訓練會將編碼器以及解碼器參數固定，並且訓練平行於解碼器的一個生成器來生成解碼器的殘差訊號(Residual Signal)，進而能夠生成較為尖銳、與真實資料較相近的輸出。



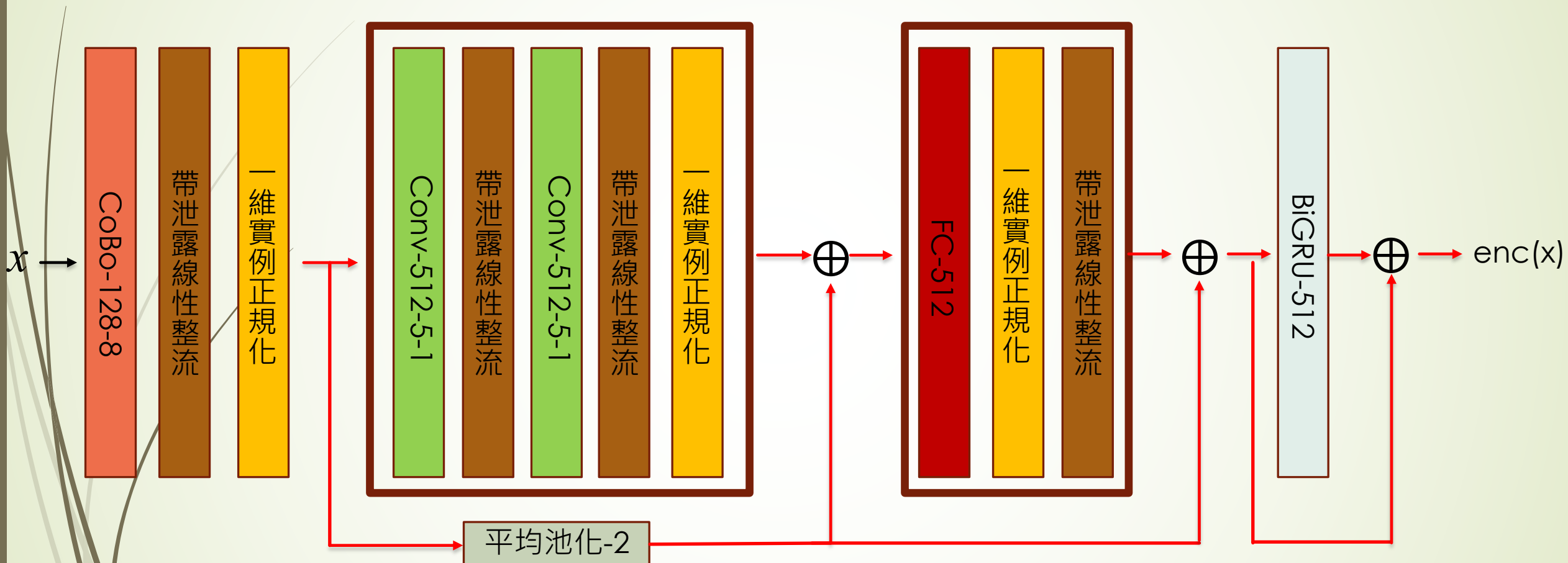
預訓練階段示意圖。在這個階段視同訓練一個語音訊號的自編碼器，但會同時在解碼器輸入語者的編號。 x 為語音訊號， y 為語者編號

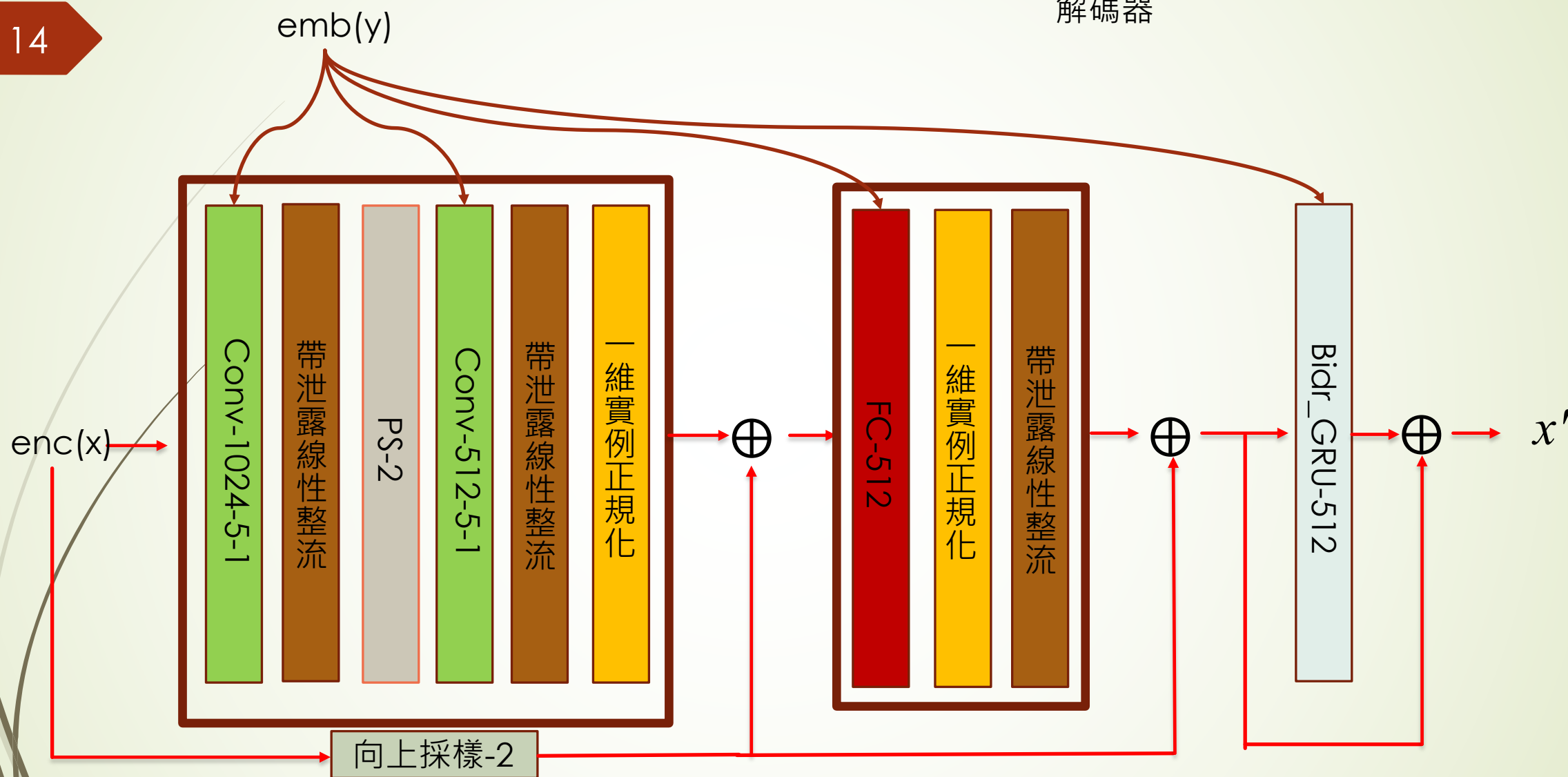
Each channel has zero mean
and unit variance

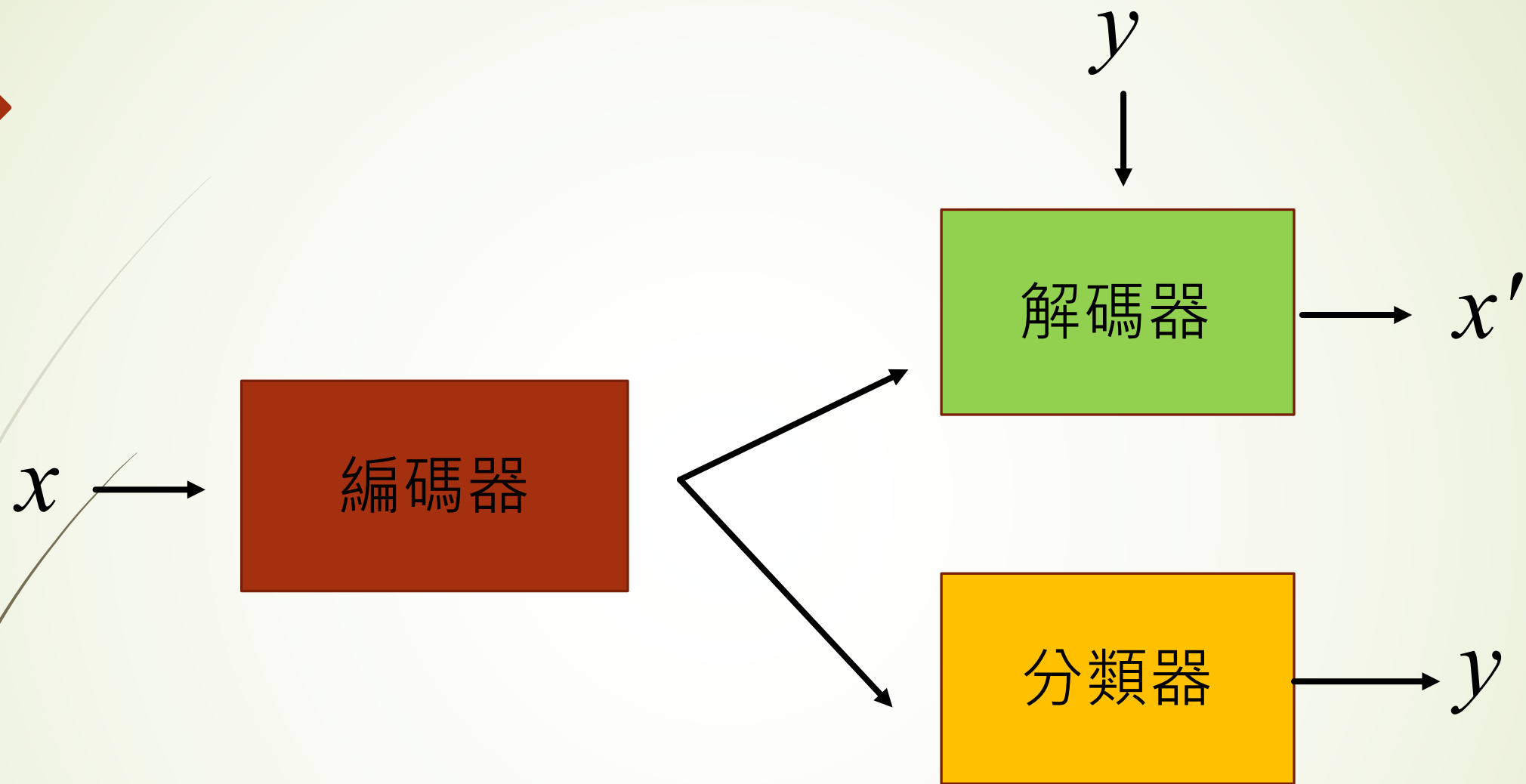
Normalize for each
channel



編碼器

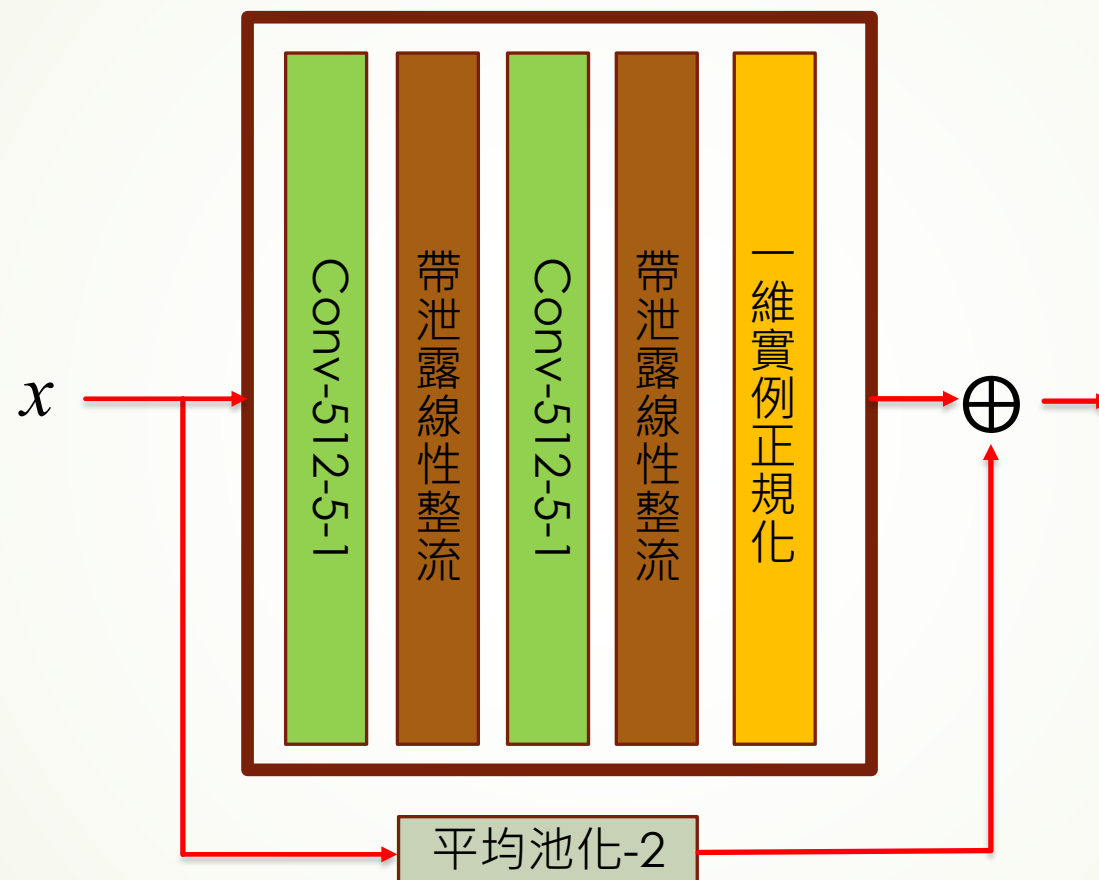


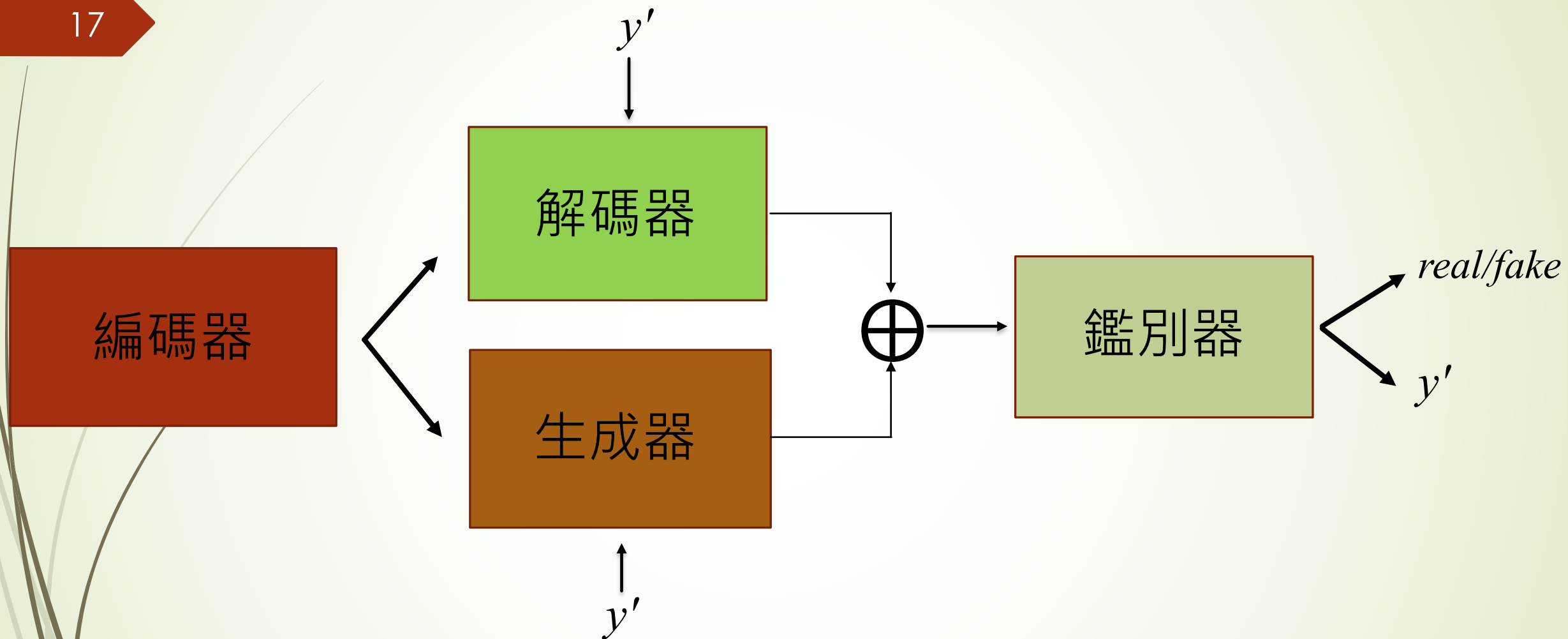




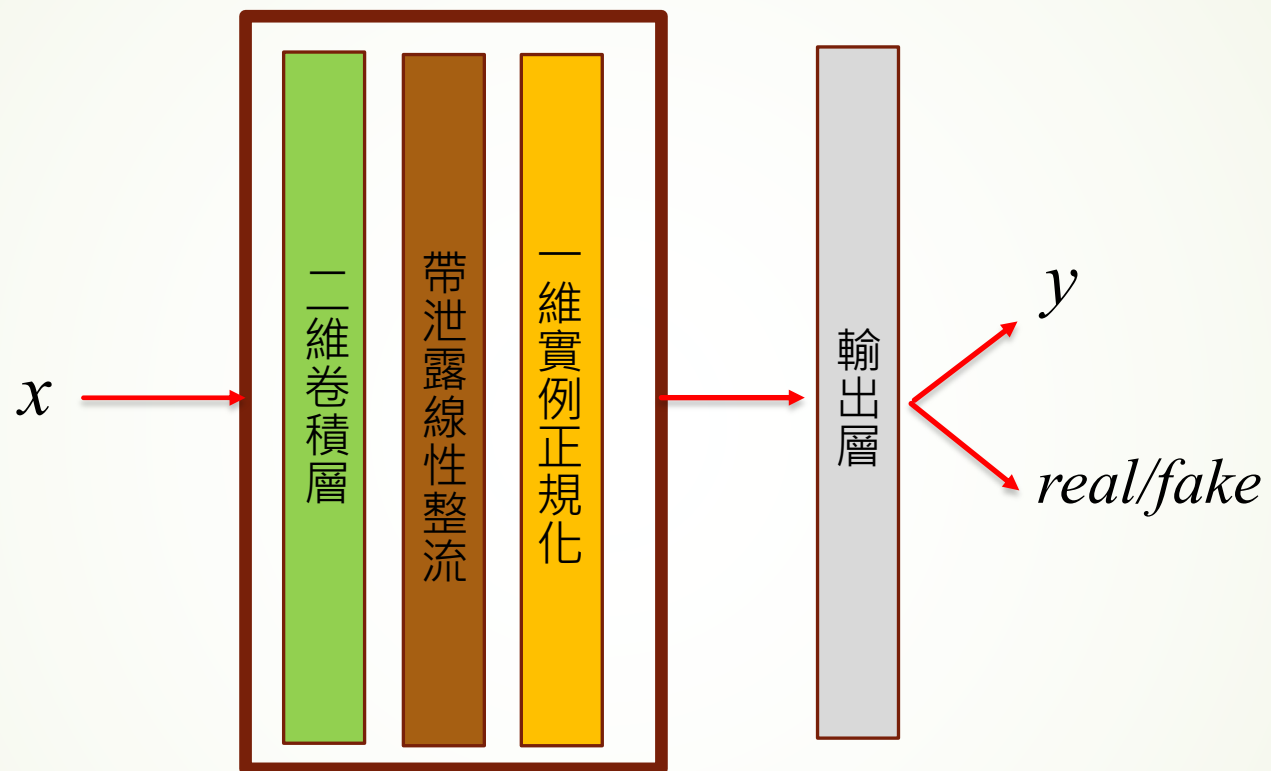
解纏特徵學習階段示意圖。

在這個階段會引入分類器來規範化(Regularize)編碼器生成的語音特徵





生成對抗網路階段。在這個階段會引入一個生成器以及一個鑑別器。透過生成器與鑑別器的交互訓練，來提升模型所生成的語音品質



DEMO

https://eric4404123.github.io/voice_conversion/