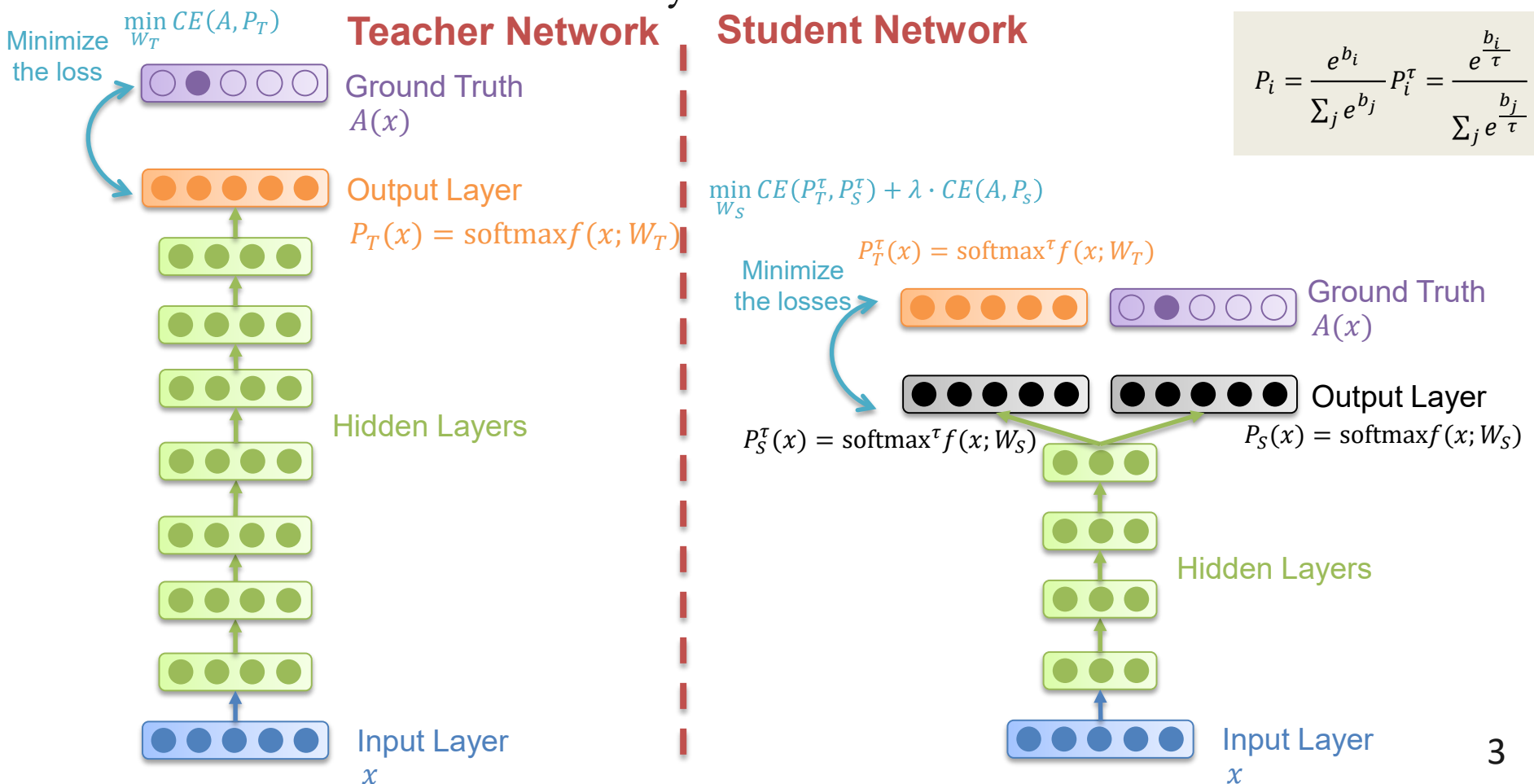# Advanced Structures

**Kuan-Yu Chen (陳冠宇)**

2018/05/10 @ NTUST

# Model Compression

- Deep networks have recently exhibited state-of-the-art performance in computer vision tasks such as image classification and object detection
  - Top-performing systems usually involve very wide and deep networks, with numerous parameters
    - time consuming
    - high memory demanding

  - Knowledge Distillation (Teacher-student network) and FitNet are representatives
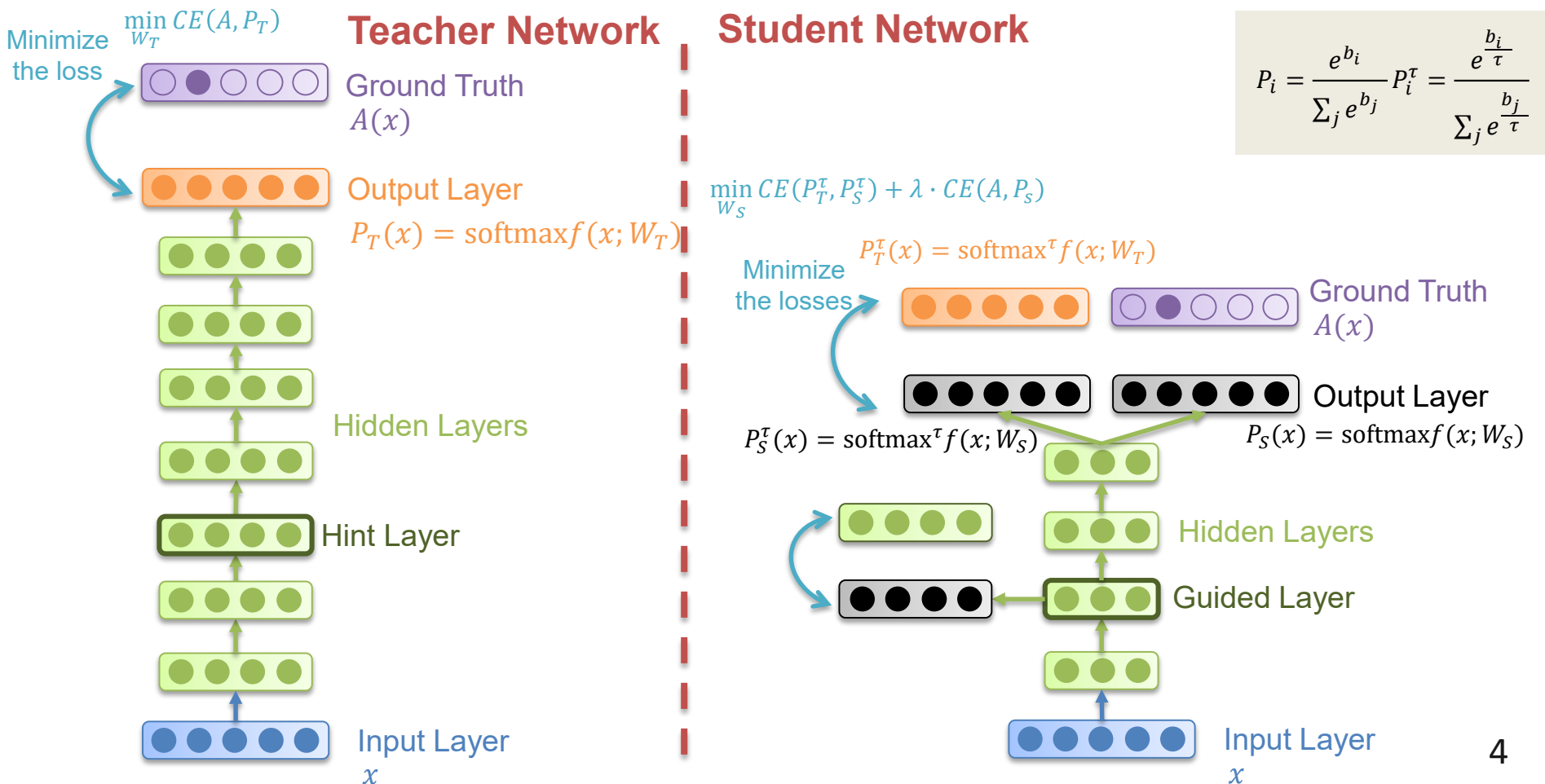
# Knowledge Distilling

- The idea is to allow the student network to capture not only the information provided by the **true labels**, but also the finer structure learned by the **teacher network**
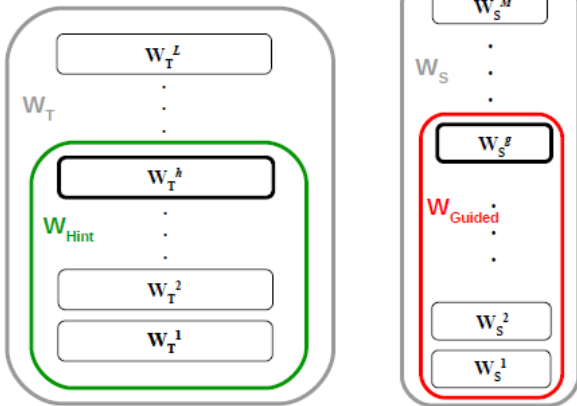
$$P_i = \frac{e^{b_i}}{\sum_j e^{b_j}} \quad P_i^\tau = \frac{e^{\frac{b_i}{\tau}}}{\sum_j e^{\frac{b_j}{\tau}}}$$

**Teacher Network**

$$\min_{W_T} CE(A, P_T)$$

Minimize the loss

Ground Truth
$A(x)$

Output Layer
$$P_T(x) = \text{softmax} f(x; W_T)$$

Hidden Layers

Input Layer
$x$

**Student Network**

$$\min_{W_S} CE(P_T^\tau, P_S^\tau) + \lambda \cdot CE(A, P_S)$$

$$P_T^\tau(x) = \text{softmax}^\tau f(x; W_T)$$

Minimize the losses

Ground Truth
$A(x)$

Output Layer

$$P_S^\tau(x) = \text{softmax}^\tau f(x; W_S)$$

$$P_S(x) = \text{softmax} f(x; W_S)$$

Hidden Layers

Input Layer
$x$

3

# FitNets.

- The FitNet is trained in a stage-wise fashion
  - The core idea is that layer-wise information should also be obtained

# FitNets..
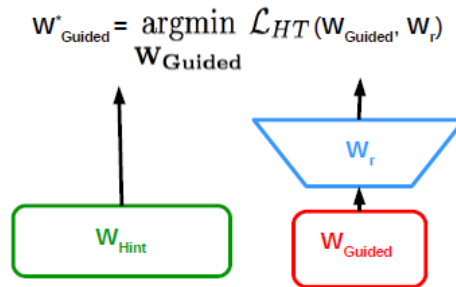


$$w^*_{Guided} = \underset{W_{Guided}}{\operatorname{argmin}} \, \mathcal{L}_{HT}(w_{Guided}, w_r)$$

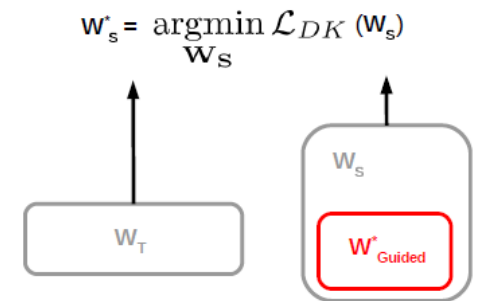$$w^*_s = \underset{W_S}{\operatorname{argmin}} \, \mathcal{L}_{DK}(w_s)$$

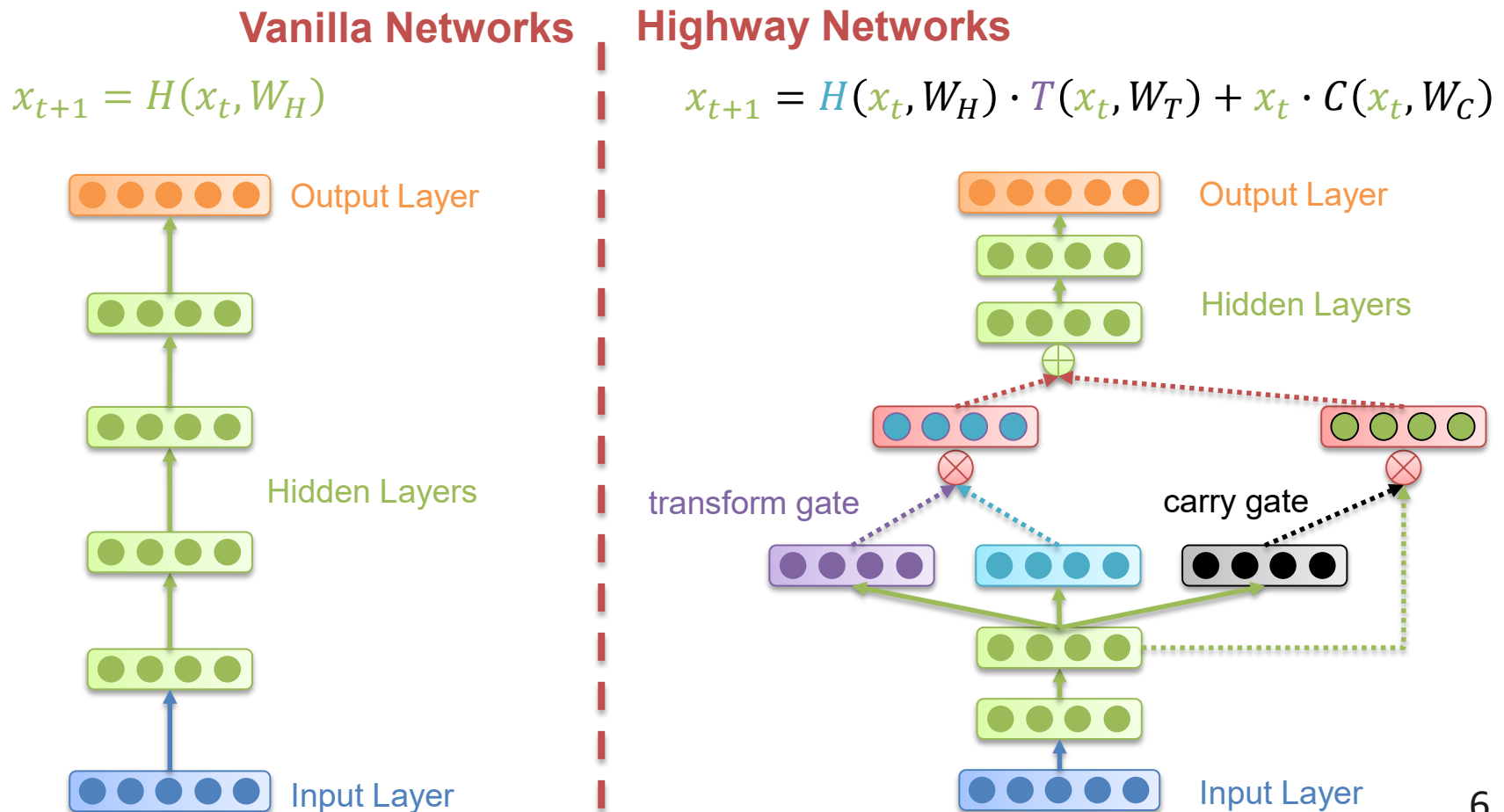(a) Teacher and Student Networks     (b) Hints Training     (c) Knowledge Distillation

1. Starting from a trained teacher network and a randomly initialized FitNet

2. Adding a regressor $W_r$ on top of the FitNet guided layer and train the FitNet parameters $W_{Guided}$

3. Based on the pre-trained parameters $W_{Guided}$, we train the parameters of whole FitNet, $W_S$

5

# **Highway Networks.**

- Highway networks allow unimpeded information flow across several layers on information highways
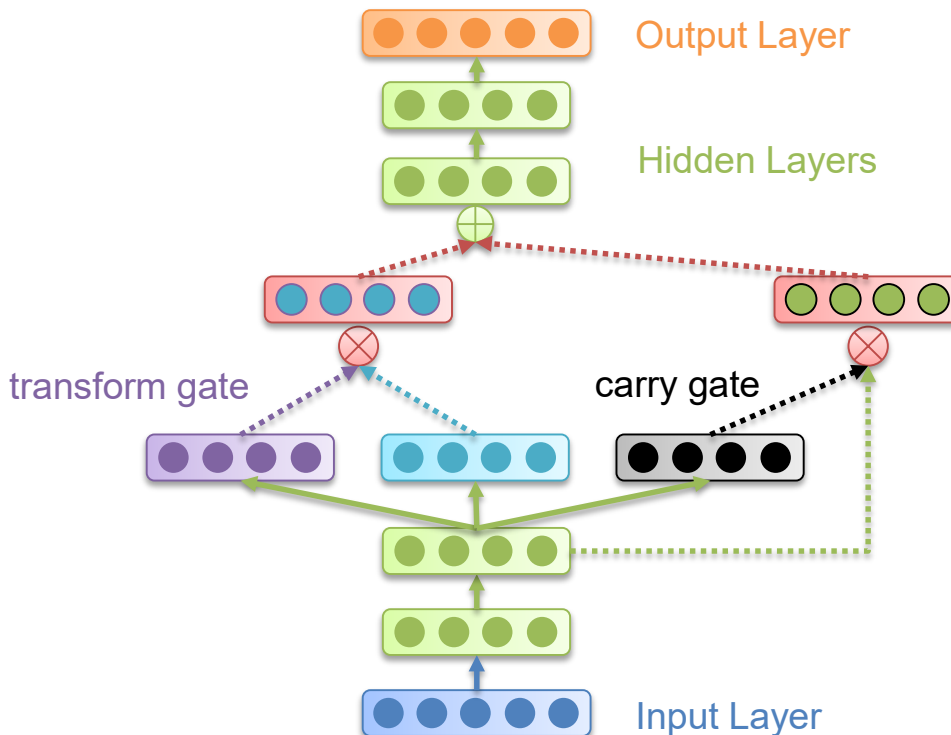
**Vanilla Networks**

$$x_{t+1} = H(x_t, W_H)$$

**Highway Networks**

$$x_{t+1} = H(x_t, W_H) \cdot T(x_t, W_T) + x_t \cdot C(x_t, W_C)$$

Output Layer

Hidden Layers

Input Layer

Output Layer

Hidden Layers

transform gate

carry gate

Input Layer

6

# Highway Networks..

- A simplified variant is to set carry gate equal to one minus transform gate
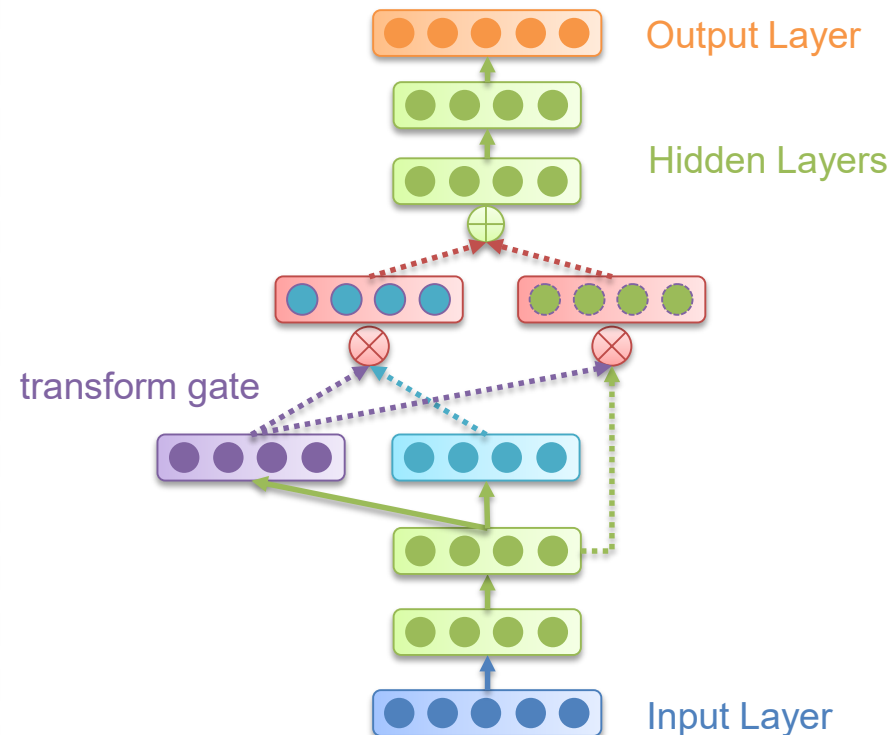
**Highway Networks**

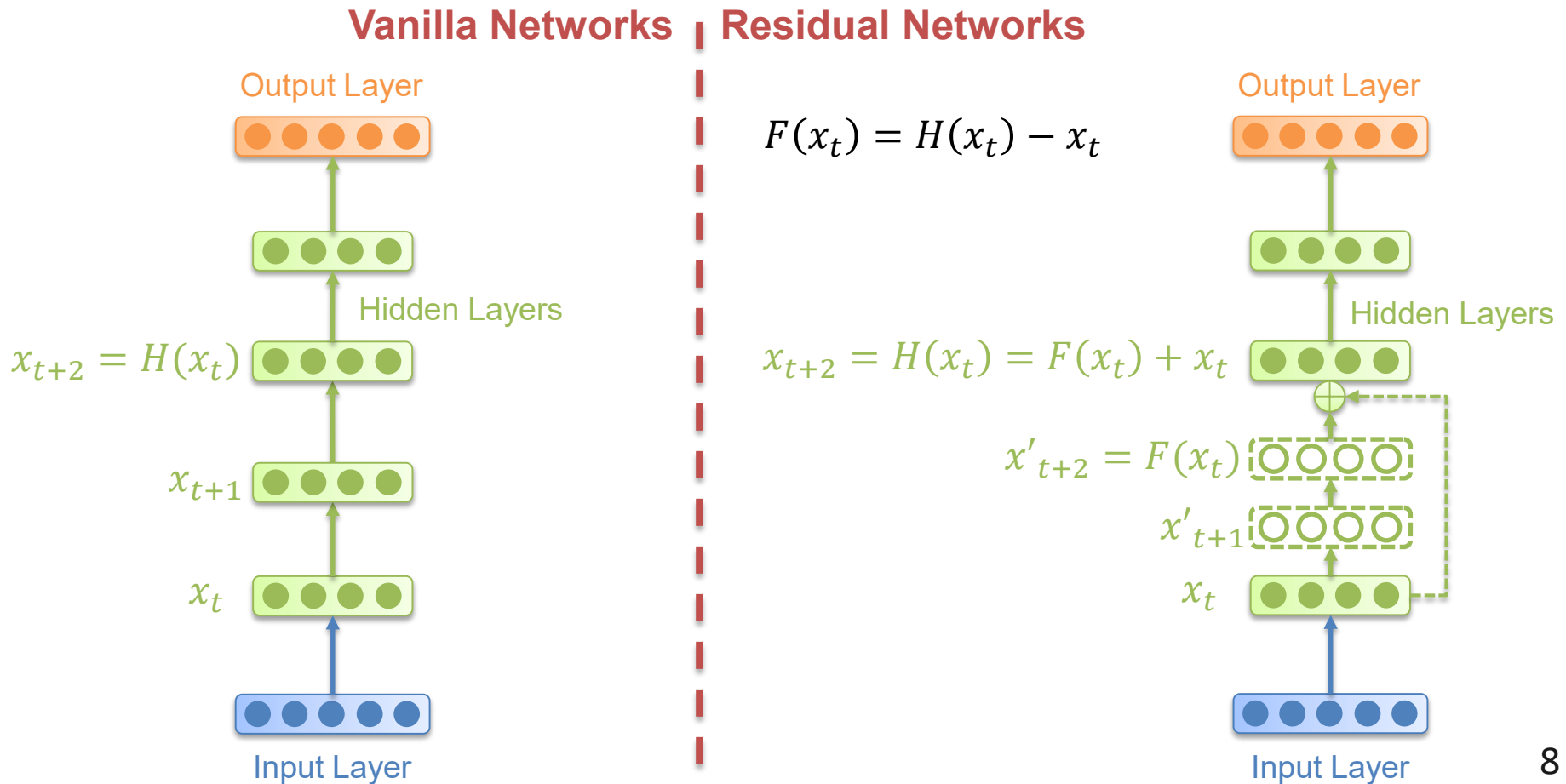$$x_{t+1} = H(x_t, W_H) \cdot T(x_t, W_T) + x_t \cdot C(x_t, W_C)$$

**Simplified Highway Networks**

$$x_{t+1} = H(x_t, W_H) \cdot T(x_t, W_T)$$
$$+ x_t \cdot [1 - T(x_t, W_T)]$$

# Residual Networks

- ResNet hypothesizes that it is easier to optimize the **residual mapping** than to optimize the original, unreferenced mapping
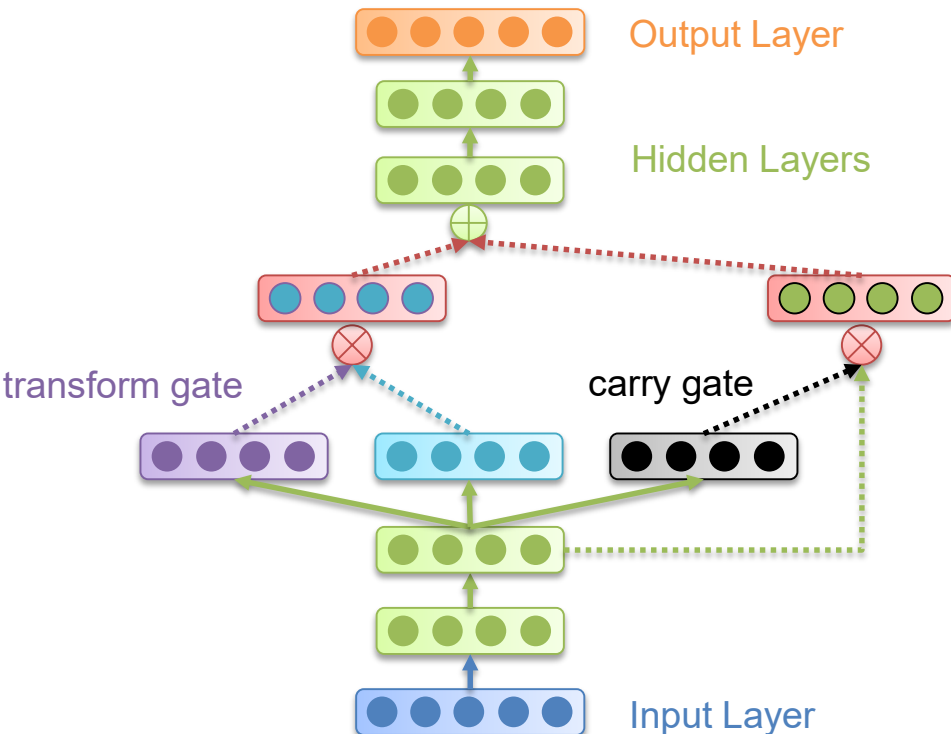


**Vanilla Networks** | **Residual Networks**

Output Layer

$$F(x_t) = H(x_t) - x_t$$

Hidden Layers

$x_{t+2} = H(x_t)$

$x_{t+2} = H(x_t) = F(x_t) + x_t$

$x'_{t+2} = F(x_t)$

$x_{t+1}$

$x'_{t+1}$

$x_t$

$x_t$

Input Layer

Output Layer

Hidden Layers

Input Layer

# Highway vs. ResNet

- ResNet is a short name for Residual Network
  - ResNet usually refers to the classic CNN-based residual learning
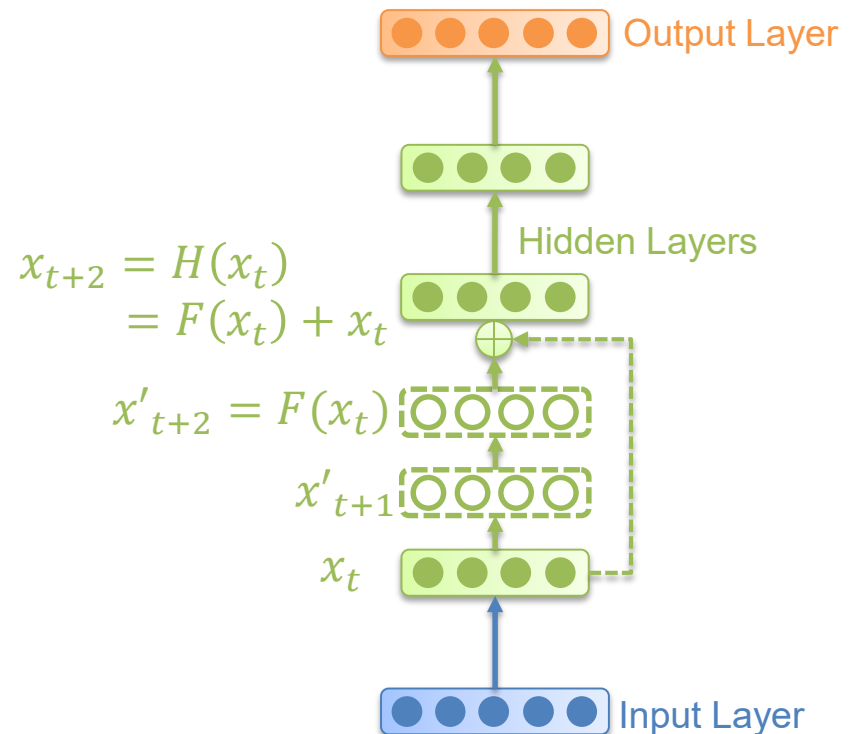- When $T(\cdot) = C(\cdot) = 1$, ResNet $\sim$ Highway

**Highway Networks** | **Residual Networks**

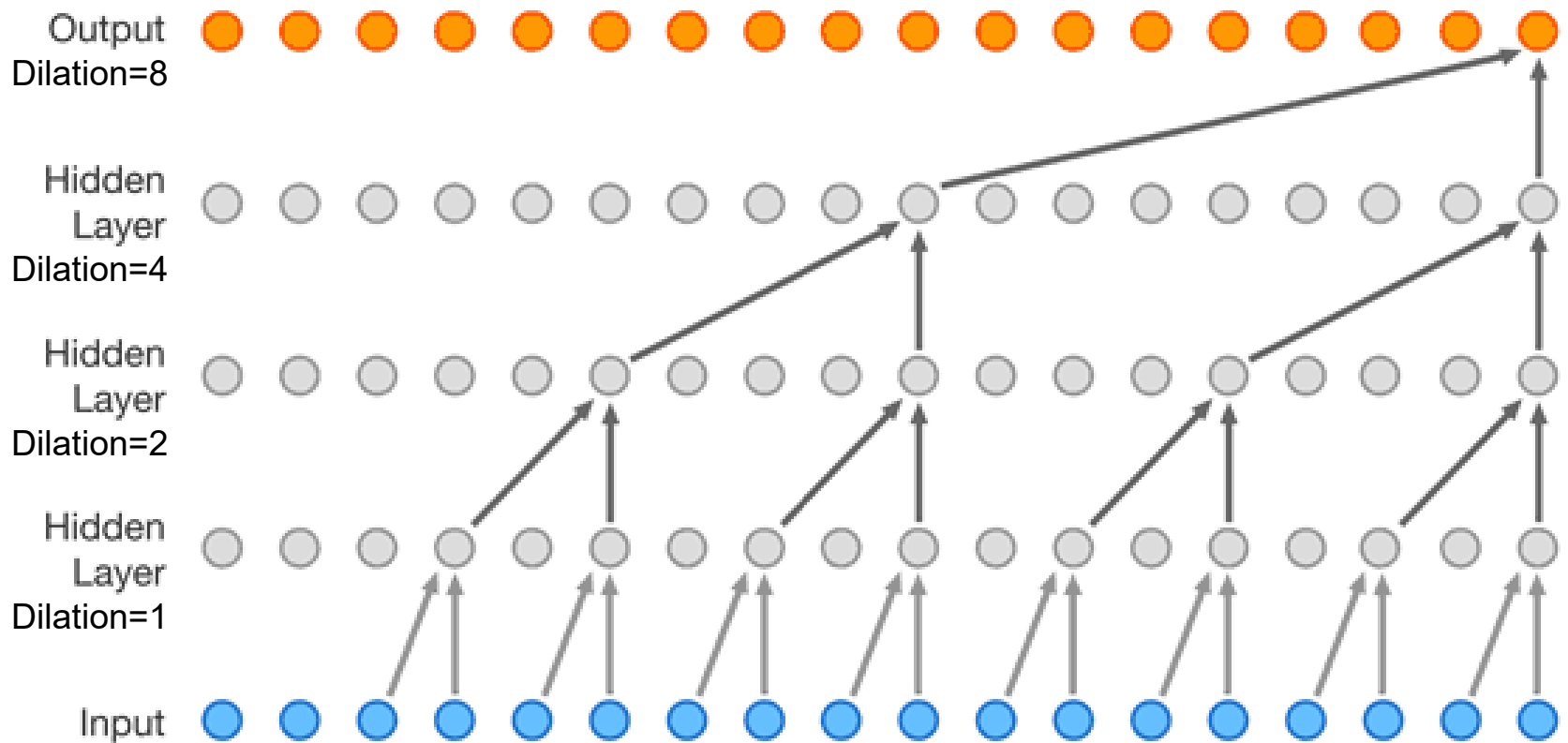$$x_{t+1} = H(x_t, W_H) \cdot T(x_t, W_T) + x_t \cdot C(x_t, W_C)$$

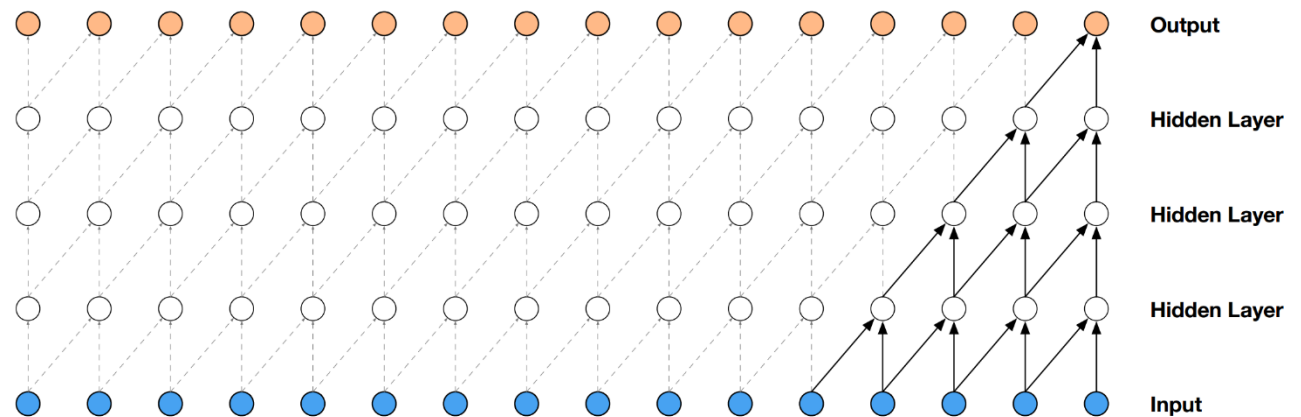$$x_{t+2} = F(x_t) + x_t \qquad F(x_t) = H(x_t) - x_t$$
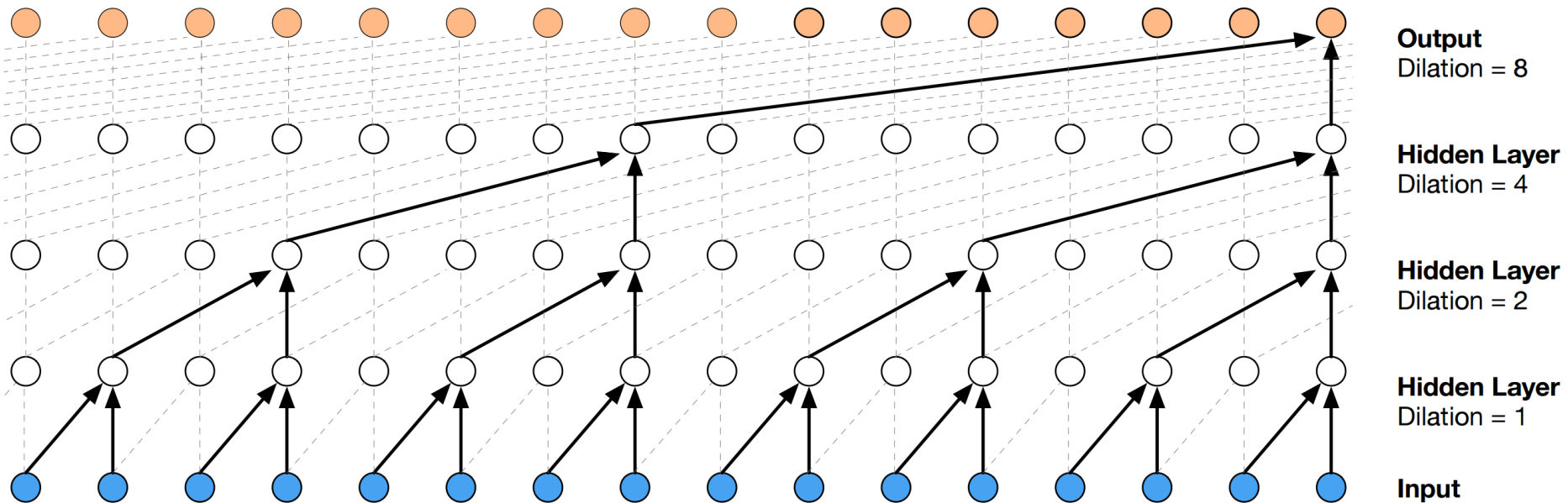
# **Dilated Convolutions.**
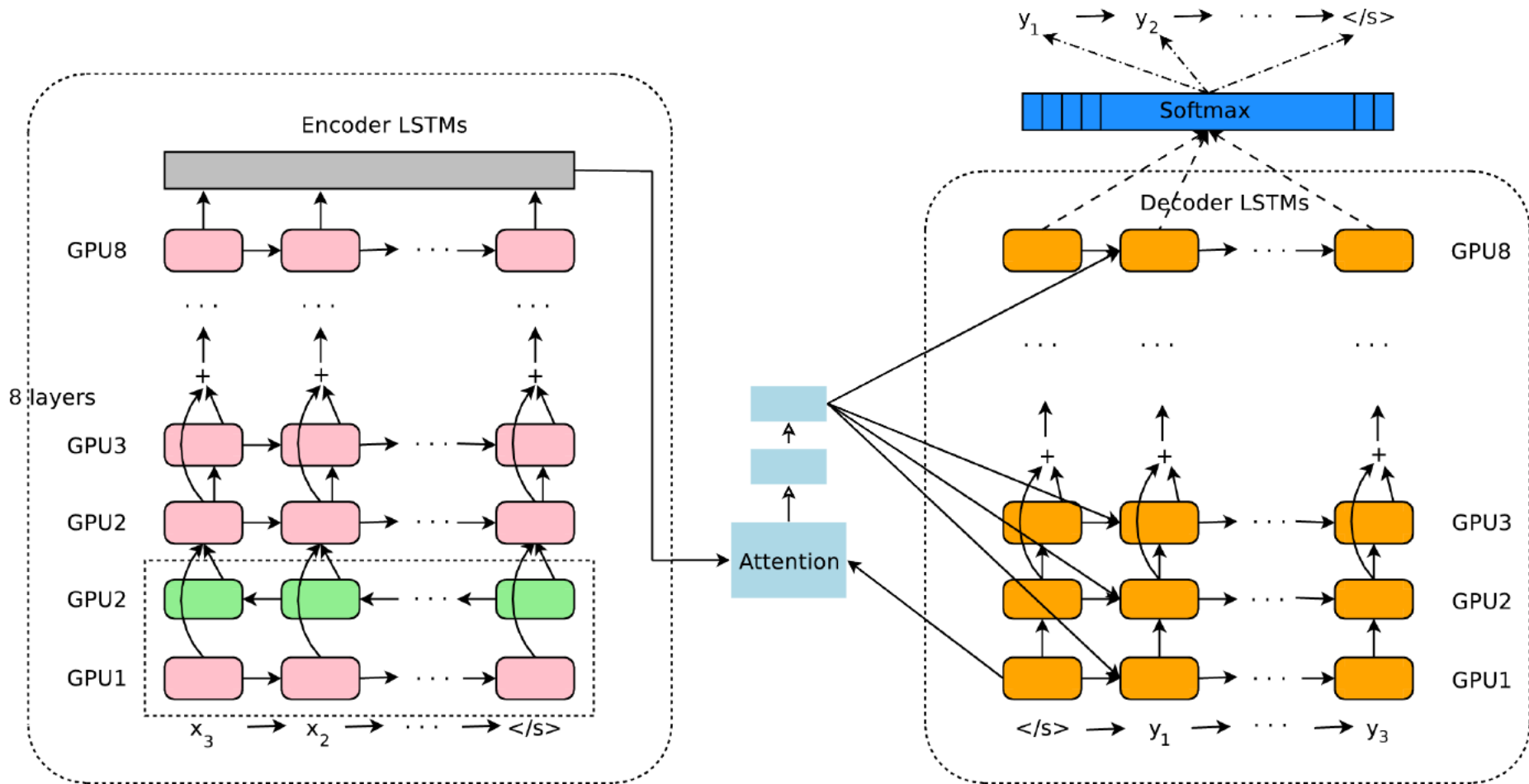
# Dilated Convolutions..

- Classic CNNs

- Dilated Convolutions

# Google's Neural Machine Translation.

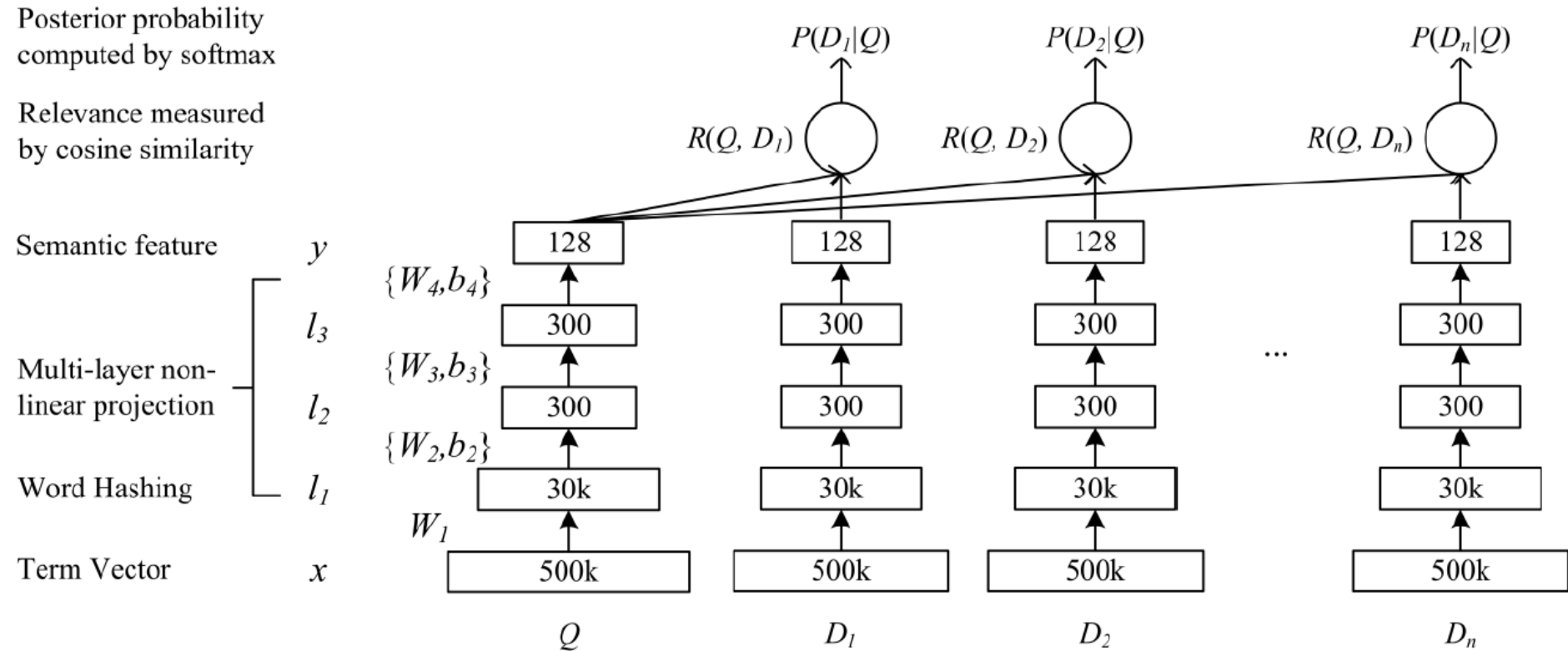- A conventional encoder-decoder architecture with attention

# Google's Neural Machine Translation..

- To be able to make use of multilingual data within a single system, GNMT proposes one simple modification to the input data
  - An artificial token is introduced at the beginning of the input sentence to indicate the target language the model should translate to

```
How are you? -> ¿Cómo estás?
```

```
<2es> How are you? -> ¿Cómo estás?
```

# Deep Structured Semantic Model (DSSM).

Posterior probability computed by softmax

Relevance measured by cosine similarity

$P(D_1|Q)$  $P(D_2|Q)$  $P(D_n|Q)$

$R(Q, D_1)$  $R(Q, D_2)$  $R(Q, D_n)$

Semantic feature  $y$  | 128 | 128 | 128 | 128 |

$\{W_4, b_4\}$  | 300 | 300 | 300 | 300 |

Multi-layer non-linear projection  $l_3$

$\{W_3, b_3\}$  $l_2$  | 300 | 300 | 300 | 300 |

$\{W_2, b_2\}$

Word Hashing  $l_1$  | 30k | 30k | 30k | 30k |

$W_1$

Term Vector  $x$  | 500k | 500k | 500k | 500k |

$Q$  $D_1$  $D_2$  $D_n$

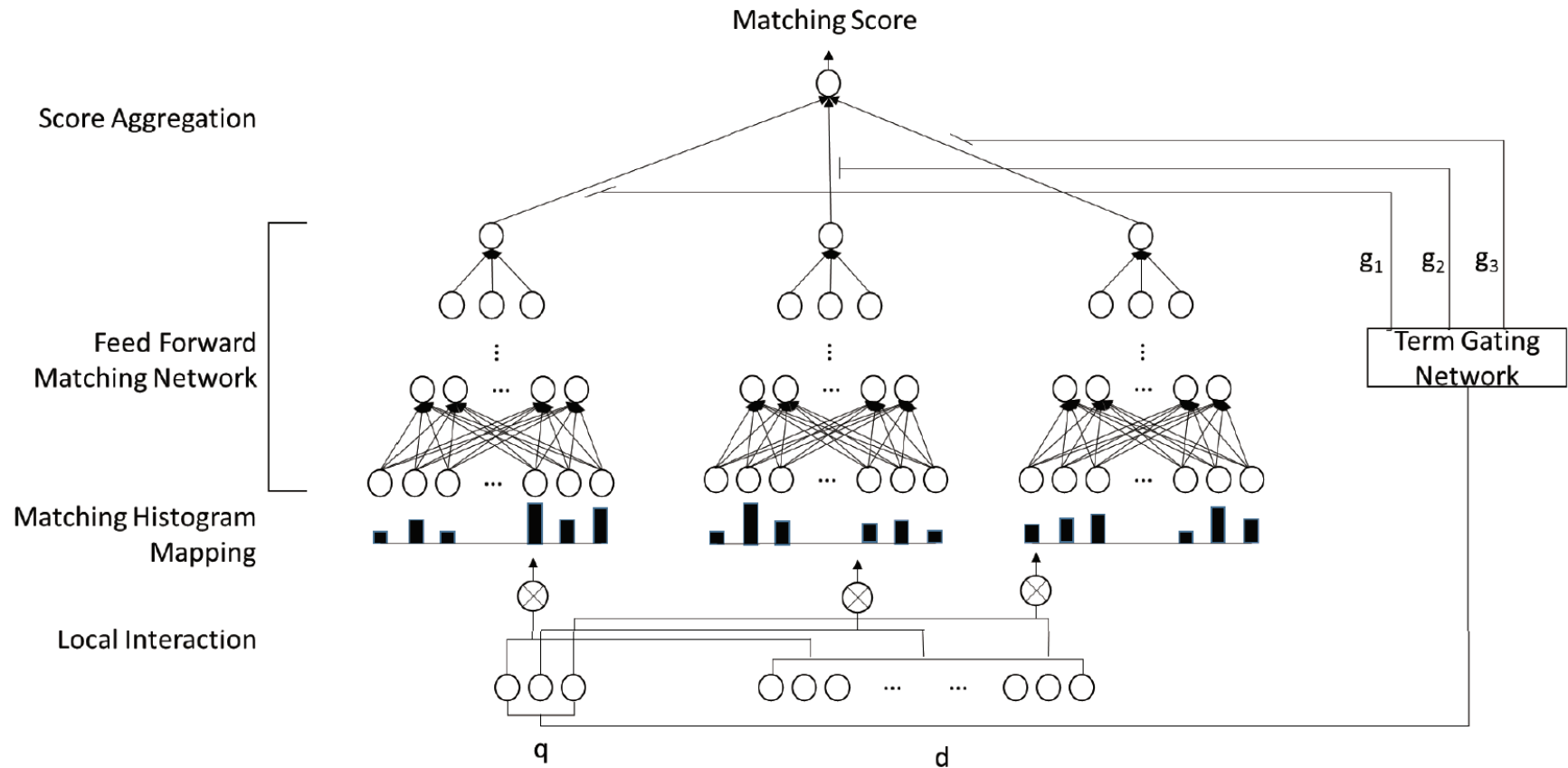| | Letter-Bigram | | Letter-Trigram | |
|---|---|---|---|---|
| Word Size | Token Size | Collision | Token Size | Collision |
| 40k | 1107 | 18 | 10306 | 2 |
| 500k | 1607 | 1192 | 30621 | 22 |

#good# => [#go, goo, ood, od#]

14

# DSSM..



$$R(q, d) = \cos(\vec{q}, \vec{d})$$

$$P(d|q) = \frac{\exp(R(q, d))}{\sum_{d'} \exp(R(q, d'))}$$

$$L = \prod_{d \in R_q} P(d|q)$$

# Deep Relevance Matching Model



Query: "car to go"
Document: "car, rent, truck, bump, injunction, runway"
Five Bins: {[-1,-0.5), [-0.5,0), [0,0.5), [0.5,1), [1,1]}
Local Interaction for "car": (1, 0.2, 0.7, 0.3, -0.1, 0.1)
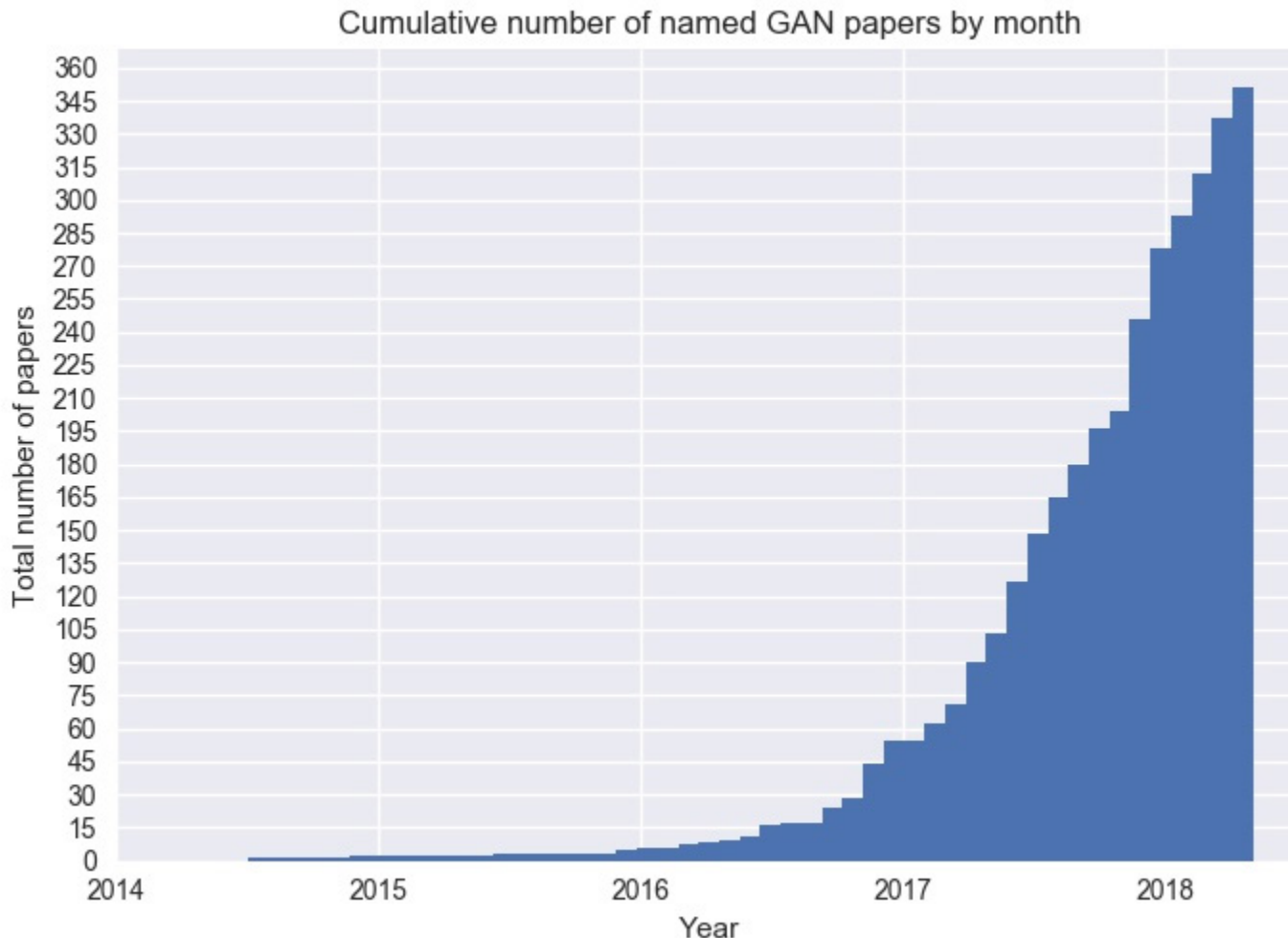Matching Histogram for "car": [0, 1, 3, 1, 1]

# Generative Adversarial Networks.

# Generative Adversarial Networks..

- https://github.com/hindupuravinash/the-gan-zoo
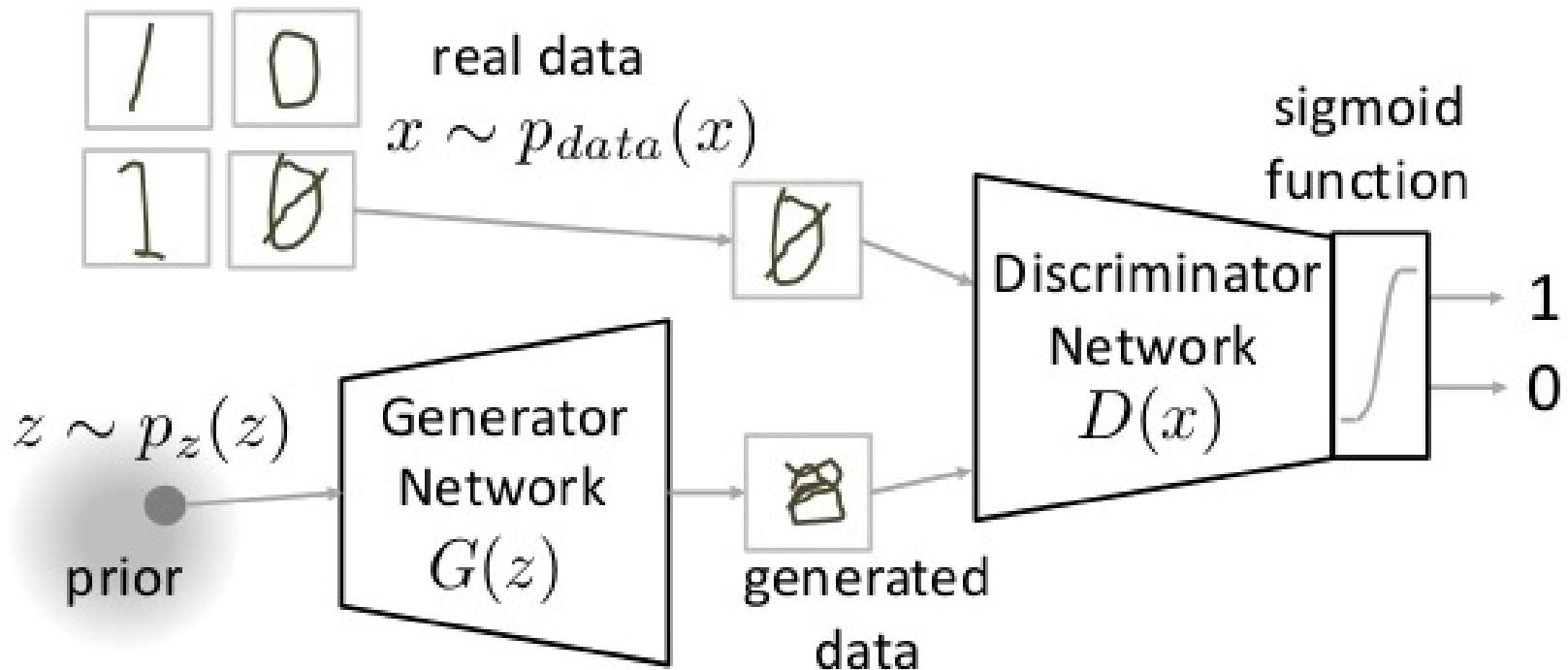


Cumulative number of named GAN papers by month
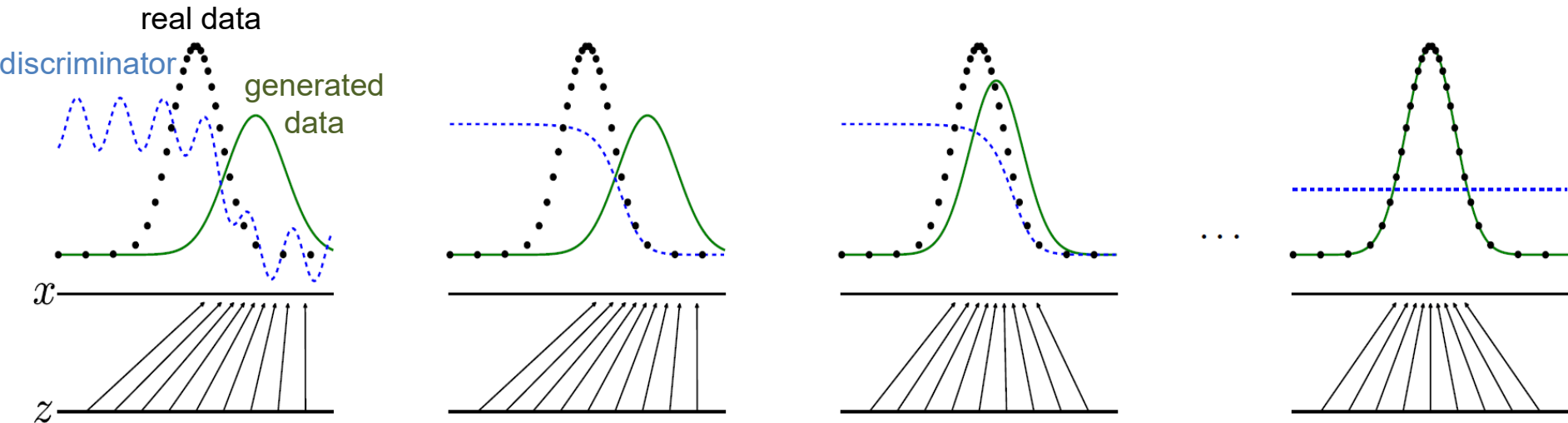
# Generative Adversarial Networks…

- Discriminator is used to criticize the results produced by generator
- The ultimate goal of generator is to cheat the discriminator, thus the generator can create potential objects

$$\min_{G} \max_{D} V(D, G) = \mathbf{E}_{x \sim p_{data}(x)}[\log D(x)] + \mathbf{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

# Generative Adversarial Networks....

$$\min_{G} \max_{D} V(D,G) = \mathbf{E}_{x \sim p_{data}(x)}[\log D(x)] + \mathbf{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

# Questions?



**kychen@mail.ntust.edu.tw**