

**Q?**

**A!**

# **FLIPDIAL**

## A GENERATIVE MODEL FOR

### **TWO-WAY**

## VISUAL DIALOGUE

**DANIELA MASSICETI, N. SIDDHARTH, PUNEET K. DOKANIA, PHILIP H.S. TORR**

UNIVERSITY OF OXFORD



CVPR, SALT LAKE CITY 2018



# TOWARD REAL-WORLD CONVERSATIONAL A.I.

*“Where is  
the bus  
stop?”*

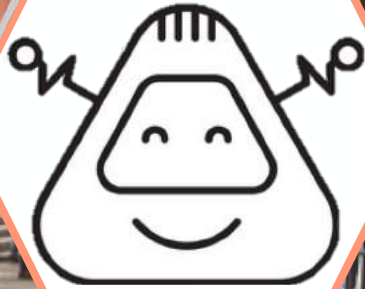




# TOWARD REAL-WORLD CONVERSATIONAL A.I.

*“Where is  
the bus  
stop?”*

*“On the  
road to  
your  
right”*





# TOWARD REAL-WORLD CONVERSATIONAL A.I.

*“Where is  
the bus  
stop?”*

*“On the  
road to  
your  
right”*



Sensible  
response



# TOWARD REAL-WORLD CONVERSATIONAL A.I.

*“Where is  
the bus  
stop?”*

*“On the  
road to  
your  
right”*



Sensible  
response

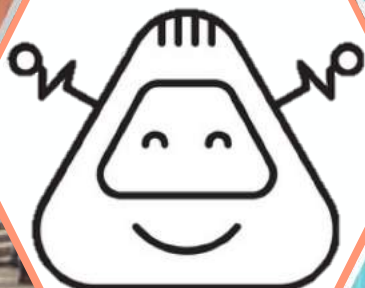
Visually  
grounded



# TOWARD REAL-WORLD CONVERSATIONAL A.I.

*“Where is  
the bus  
stop?”*

*“On the  
road to  
your  
right”*



Sensible  
response

Visually  
grounded

Natural



# TOWARD REAL-WORLD CONVERSATIONAL A.I.

*“Where is  
the bus  
stop?”*

*“On the  
road to  
your  
right”*

Sensible  
response

Correct  
grammar

Visually  
grounded

Natural





Where is  
the bus  
stop?






Where is  
the bus  
stop?

On the  
road to  
your  
right








Where is  
the bus  
stop?

On the  
road to  
your  
right

Where  
do you  
want to  
go?



A vibrant, slightly overcast scene of a busy London street, likely Piccadilly Circus. A red double-decker bus is prominent in the middle ground, displaying 'MANGO' and 'Liverpool Street 23'. The background features grand, historic buildings with domes and arches. Pedestrians are walking on the sidewalks, and a black metal railing is in the foreground. Overlaid on the left are four hexagonal callouts with a faint background pattern.

Where is  
the bus  
stop?

On the  
road to  
your  
right

Where  
do you  
want to  
go?

Home



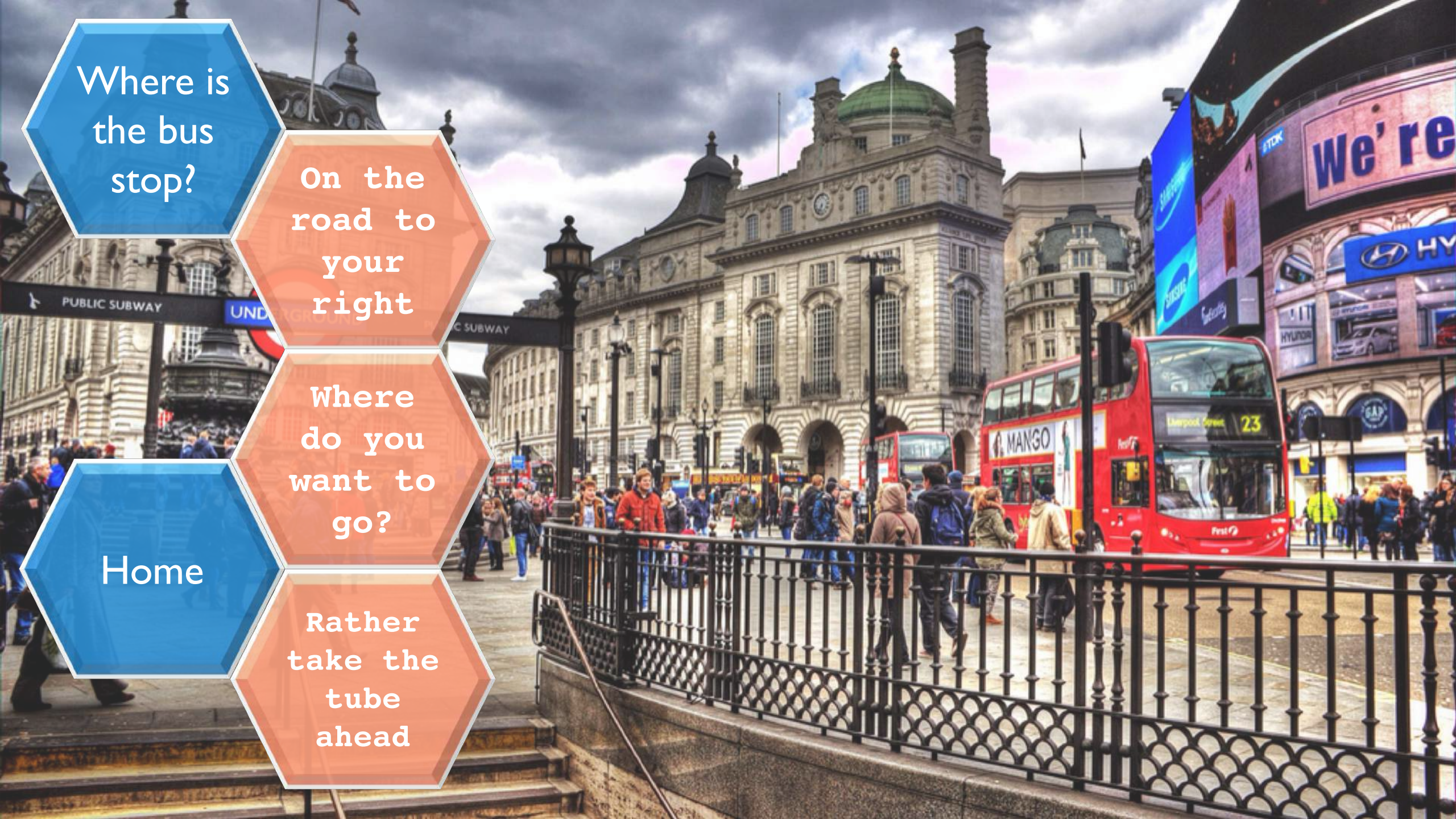
Where is  
the bus  
stop?

On the  
road to  
your  
right

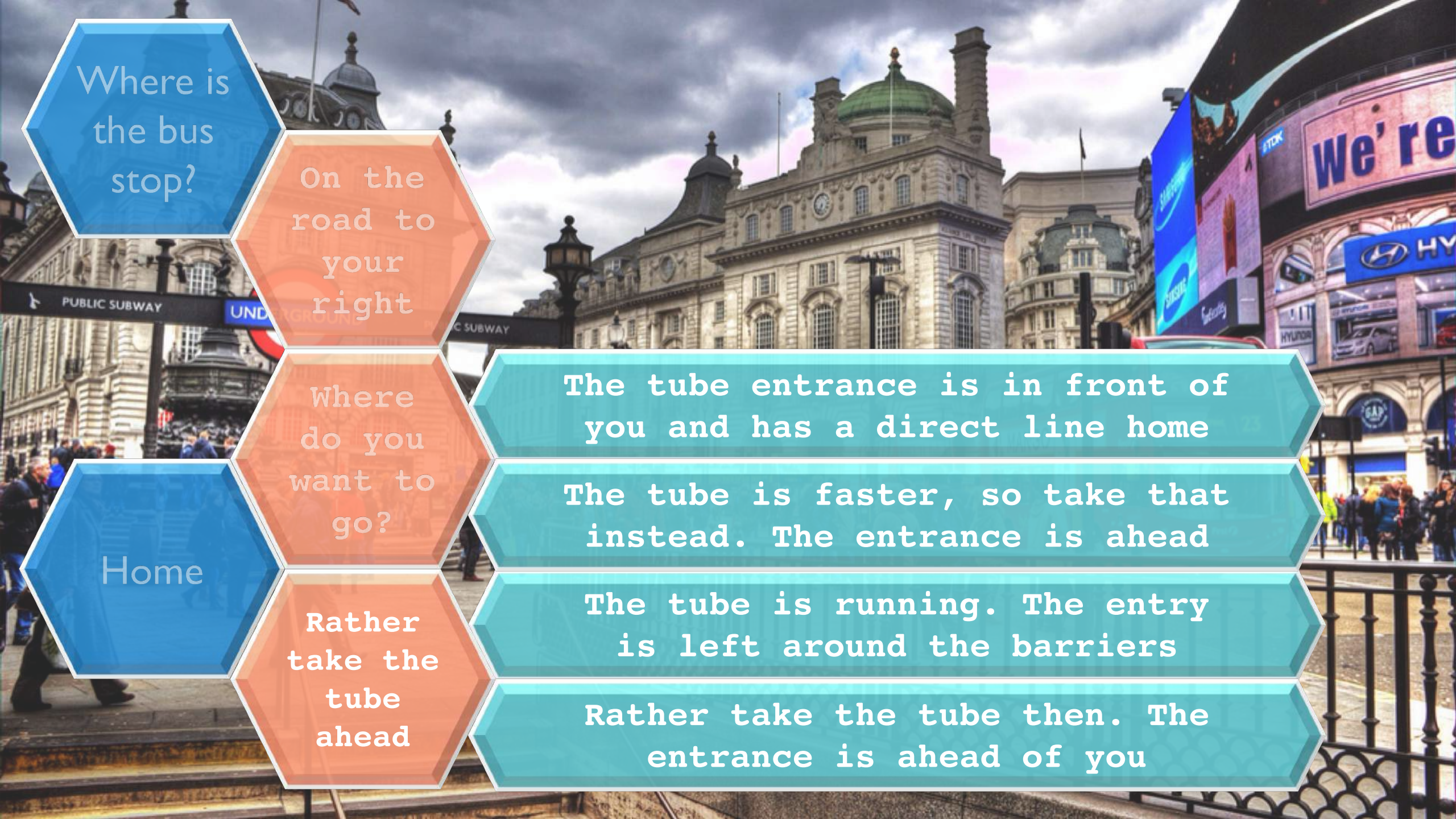
Where  
do you  
want to  
go?

Home

Rather  
take the  
tube  
ahead







Where is  
the bus  
stop?

On the  
road to  
your  
right

Where  
do you  
want to  
go?

Home

Rather  
take the  
tube  
ahead

The tube entrance is in front of  
you and has a direct line home

The tube is faster, so take that  
instead. The entrance is ahead

The tube is running. The entry  
is left around the barriers

Rather take the tube then. The  
entrance is ahead of you



## ONE-WAY VISUAL DIALOGUE (IVD)



A street with people and balloons

## TWO-WAY VISUAL DIALOGUE (2VD)



A street with people and balloons



## ONE-WAY VISUAL DIALOGUE (IVD)

How many people  
are there?

it's crowded

Where is the  
pavement?

on the left

Where are the  
balloons?

in the air

A street with people and balloons

## TWO-WAY VISUAL DIALOGUE (2VD)



A street with people and balloons



## ONE-WAY VISUAL DIALOGUE (IVD)

How many people  
are there?

it's crowded

Where is the  
pavement?

on the left

Where are the  
balloons?

in the air

A street with people and balloons

## TWO-WAY VISUAL DIALOGUE (2VD)

How many people  
are there?

it's crowded

Where is the  
pavement?

on the left

Where are the  
balloons?

in the air

A street with people and balloons



# CONDITIONAL VAE FOR IVD

Sohn et al. (2015), Kingma & Welling (2014)



image  $i$

A girl in a pink shirt eating a frosted donut

caption  $c$

|  |      |
|--|------|
| How old does the girl appear to be?        | Four |
| Does the donut have sprinkles?             | No   |
| Are there any other people with her?       | No   |
| <b>What color hair does the girl have?</b> |      |

dialogue history  $h_t^+ = \{h_{t-1}, q_t\}$

$$\log p_{\theta}(a_t \mid i, c, h_t^+)$$

Blonde



# CONDITIONAL VAE FOR IVD

Sohn et al. (2015), Kingma & Welling (2014)



image  $i$

A girl in a pink shirt eating a frosted donut

caption  $c$

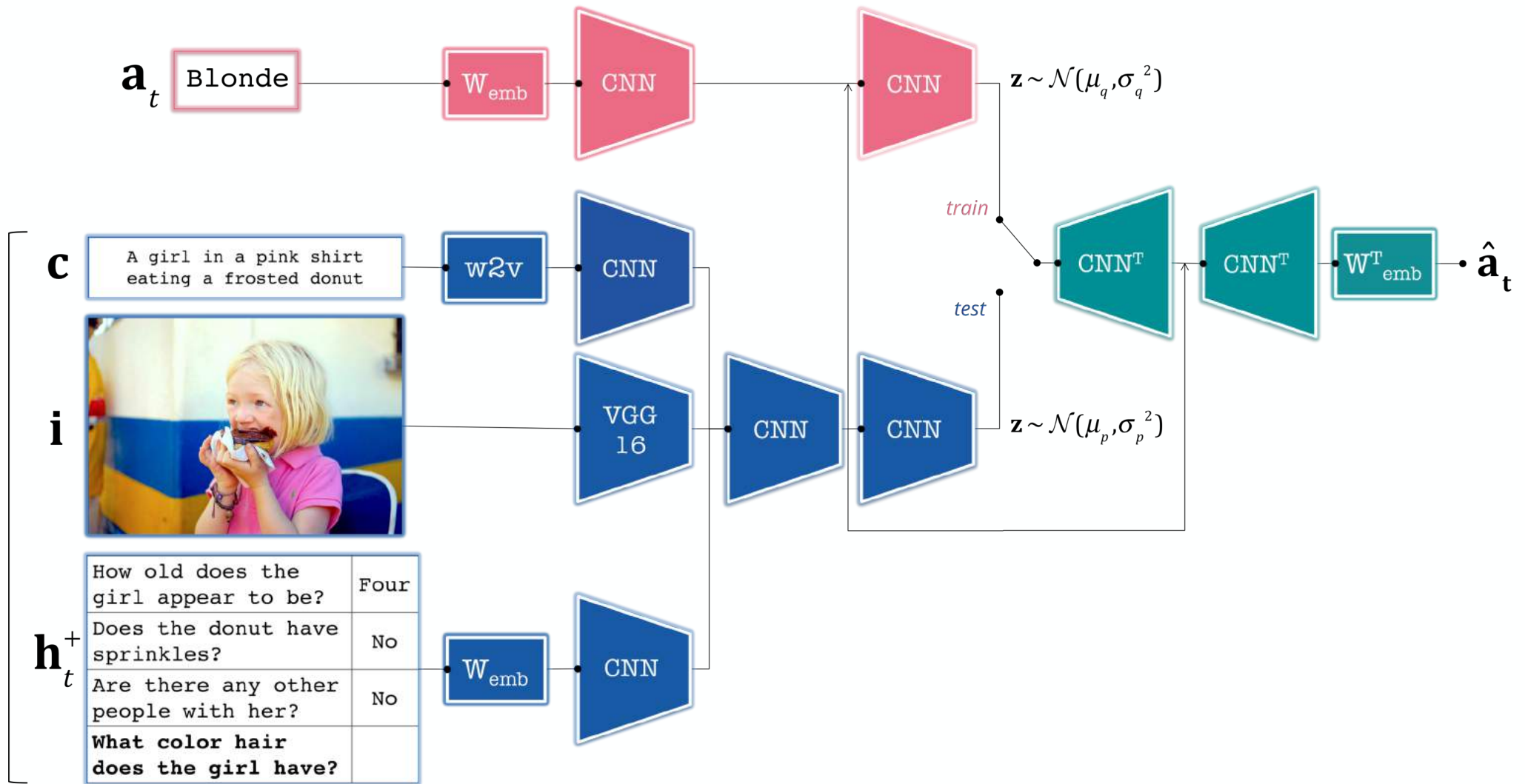
|  |      |
|--|------|
| How old does the girl appear to be?        | Four |
| Does the donut have sprinkles?             | No   |
| Are there any other people with her?       | No   |
| <b>What color hair does the girl have?</b> |      |

dialogue history  $h_t^+ = \{h_{t-1}, q_t\}$

$$\log p_\theta(a_t \mid i, c, h_t^+) \geq \mathbb{E}_{q_\phi(z \mid a_t, i, c, h_t^+)} [\log p_\theta(a_t \mid z, i, c, h_t^+)] - \mathbb{D}_{\text{KL}}(q_\phi(z \mid a_t, i, c, h_t^+) \parallel p_\theta(z \mid i, c, h_t^+))$$

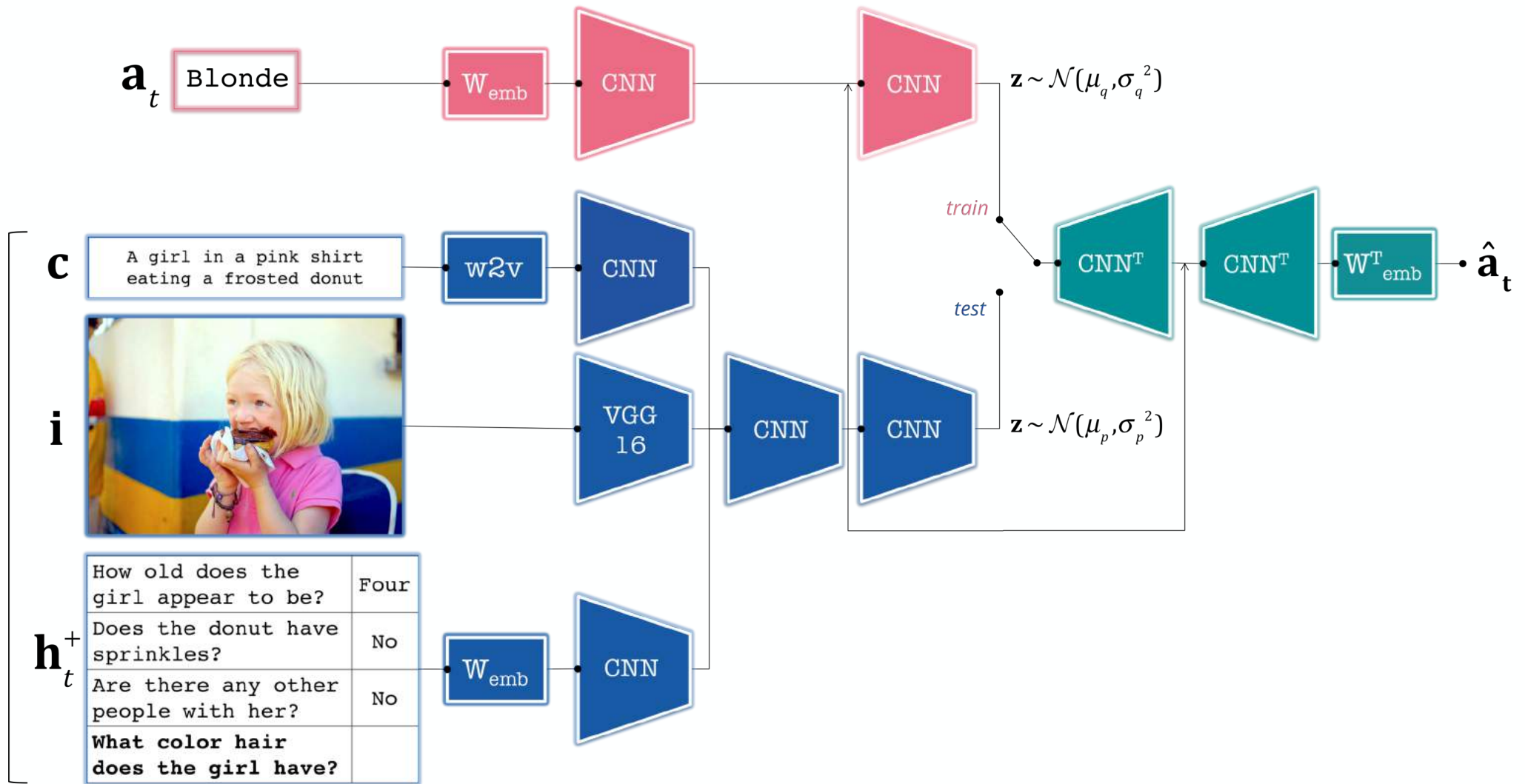
Blonde





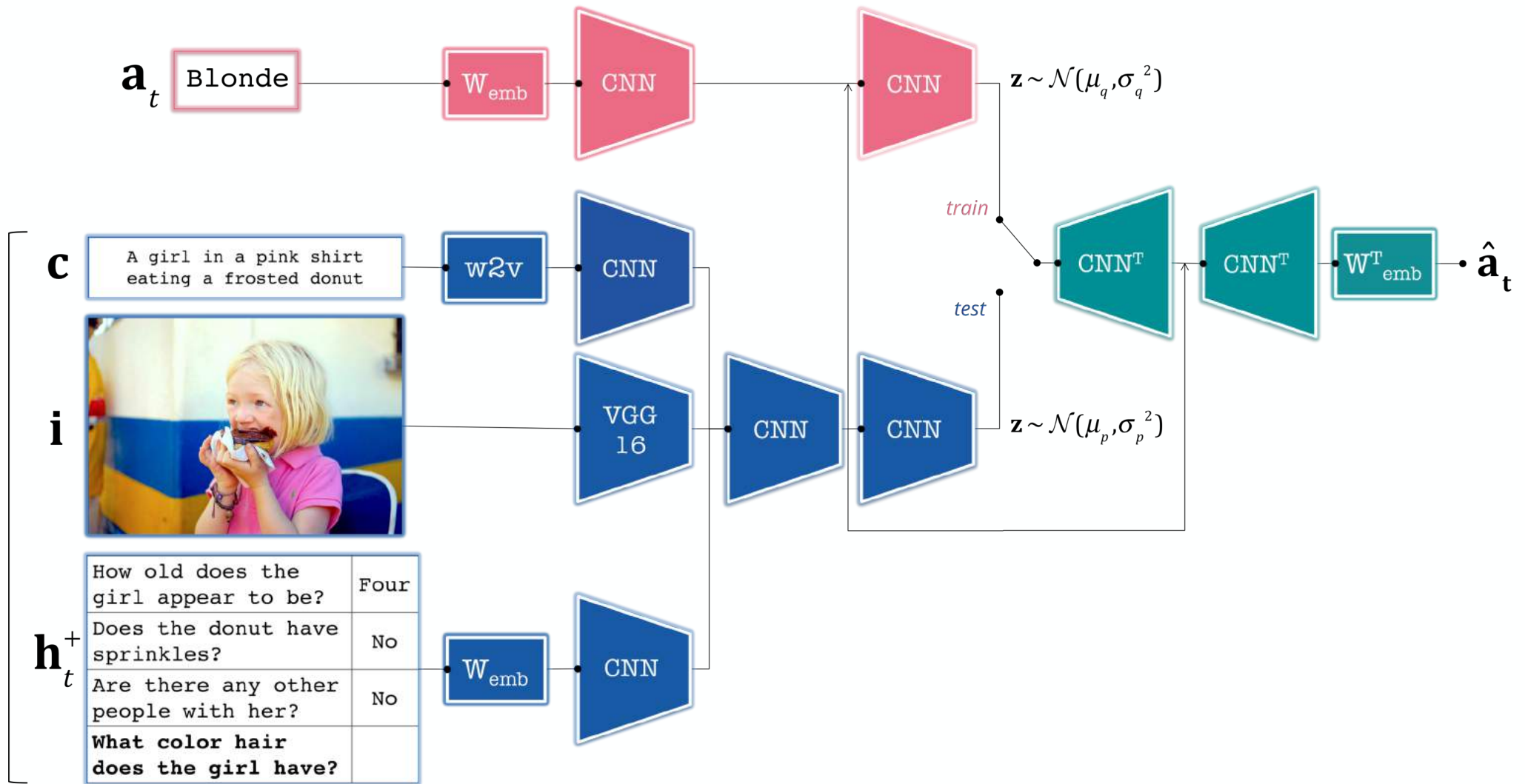
$$\log p_{\theta}(\mathbf{a}_t \mid \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+) \geq \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{a}_t, \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+)} [\log p_{\theta}(\mathbf{a}_t \mid \mathbf{z}, \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+)] \\ - \mathbb{D}_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{a}_t, \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+))$$





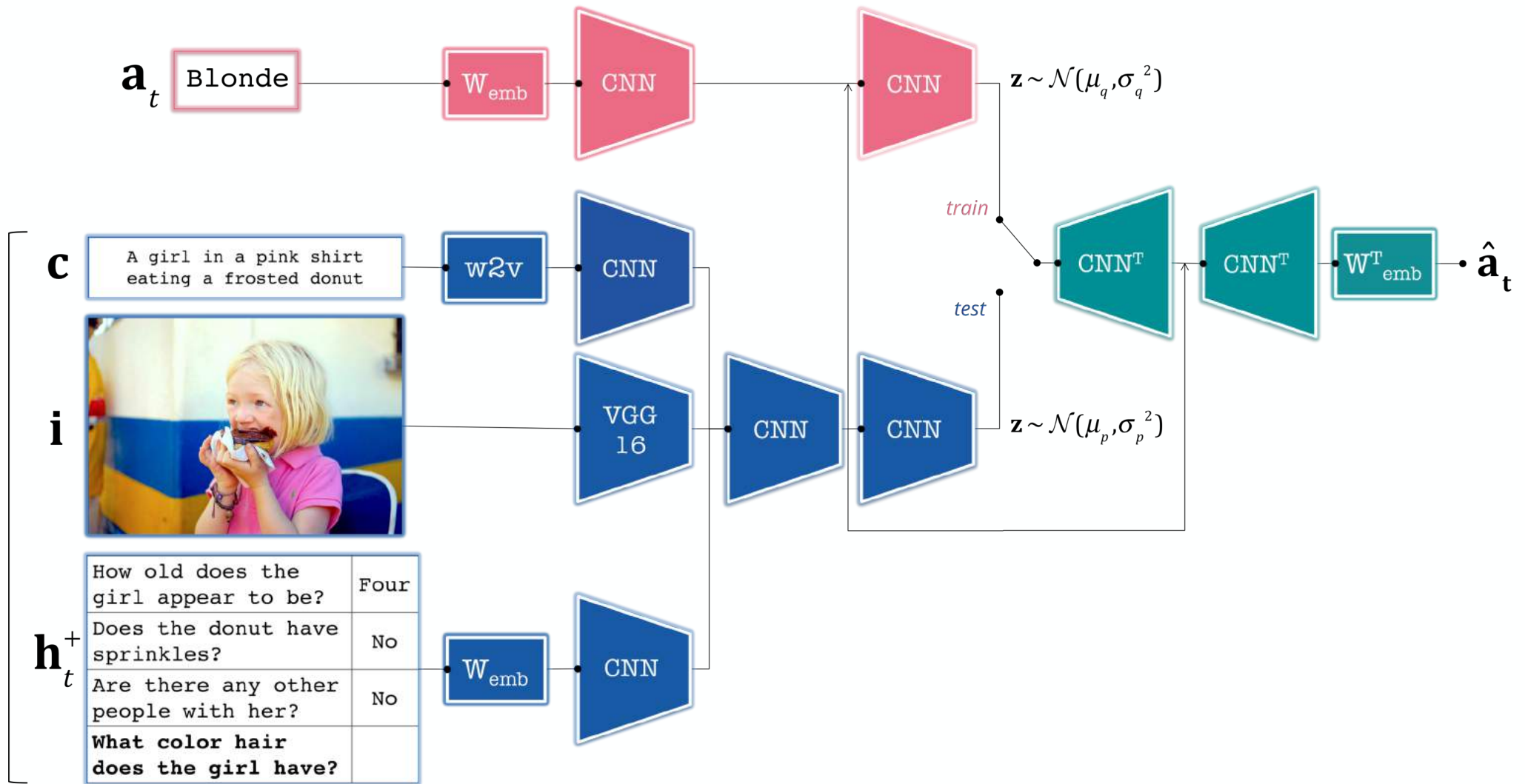
$$\log p_{\theta}(\mathbf{a}_t \mid \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+) \geq \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{a}_t, \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+)} [\log p_{\theta}(\mathbf{a}_t \mid \mathbf{z}, \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+)] \\ - \mathbb{D}_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{a}_t, \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+))$$





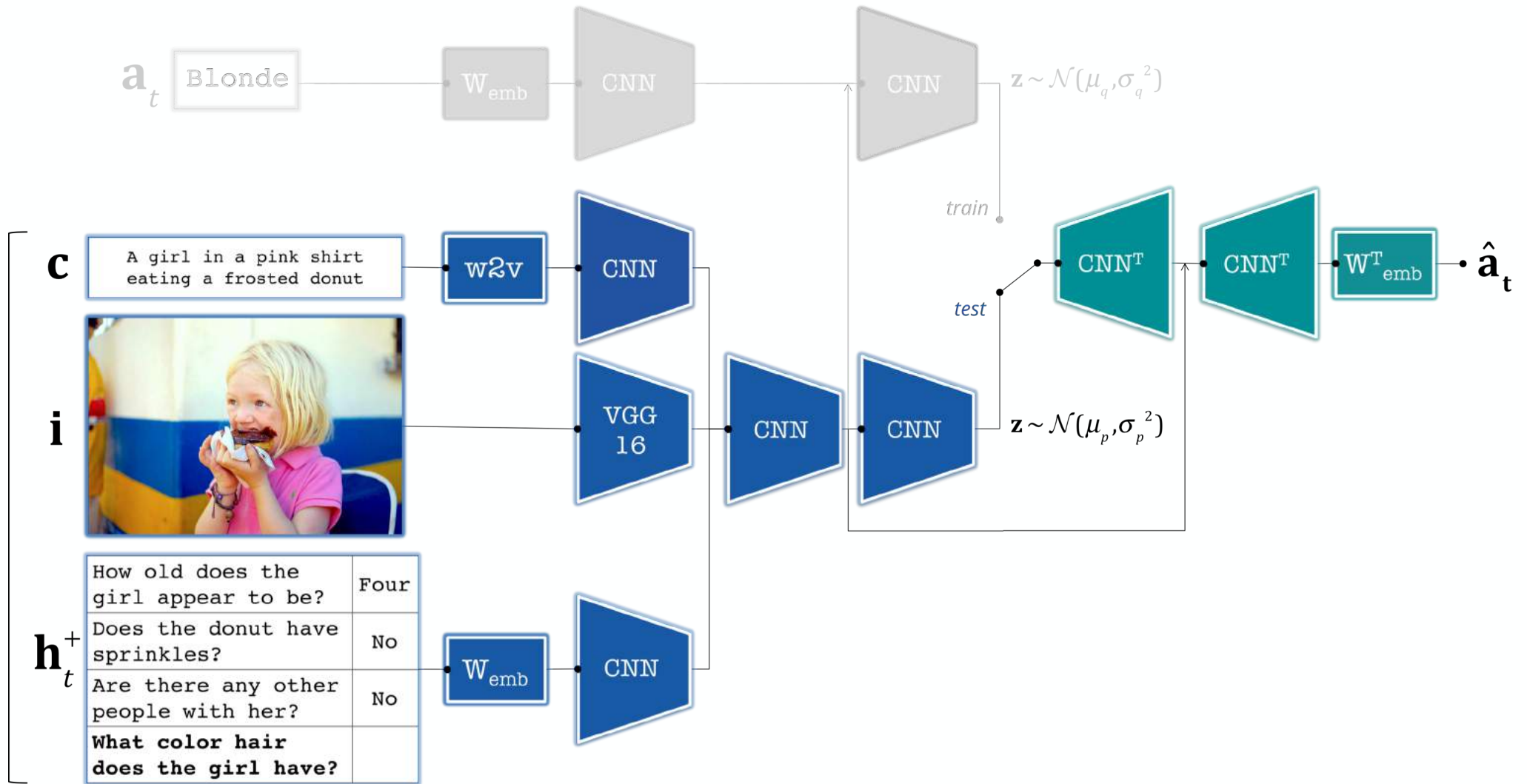
$$\log p_{\theta}(\mathbf{a}_t \mid \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+) \geq \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{a}_t, \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+)} [\log p_{\theta}(\mathbf{a}_t \mid \mathbf{z}, \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+)] - \mathbb{D}_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{a}_t, \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{i}, \mathbf{c}, \mathbf{h}_t^+))$$





$$\log p_{\theta}(a_t \mid i, c, h_t^+) \geq \mathbb{E}_{q_{\phi}(z \mid a_t, i, c, h_t^+)} [\log p_{\theta}(a_t \mid z, i, c, h_t^+)] - \mathbb{D}_{\text{KL}}(q_{\phi}(z \mid a_t, i, c, h_t^+) \parallel p_{\theta}(z \mid i, c, h_t^+))$$

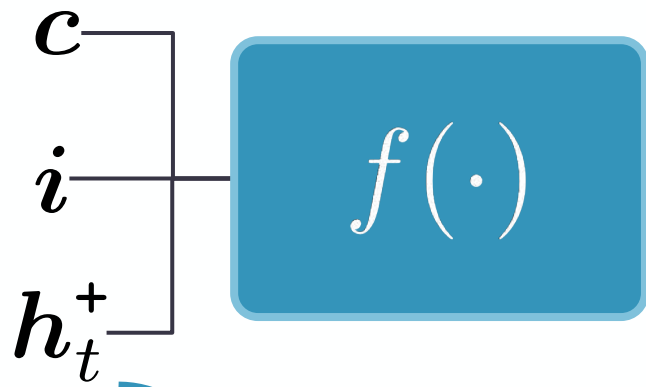




$$\log p_{\theta}(a_t \mid i, c, h_t^+) \geq \mathbb{E}_{q_{\phi}(z \mid a_t, i, c, h_t^+)} [\log p_{\theta}(a_t \mid z, i, c, h_t^+)] - \mathbb{D}_{\text{KL}}(q_{\phi}(z \mid a_t, i, c, h_t^+) \parallel p_{\theta}(z \mid i, c, h_t^+))$$



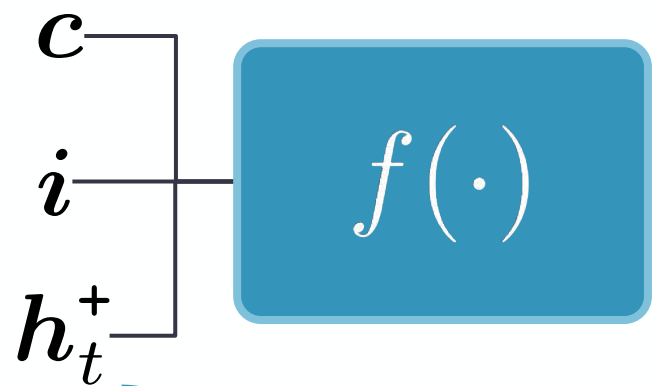
# IVD EVALUATION



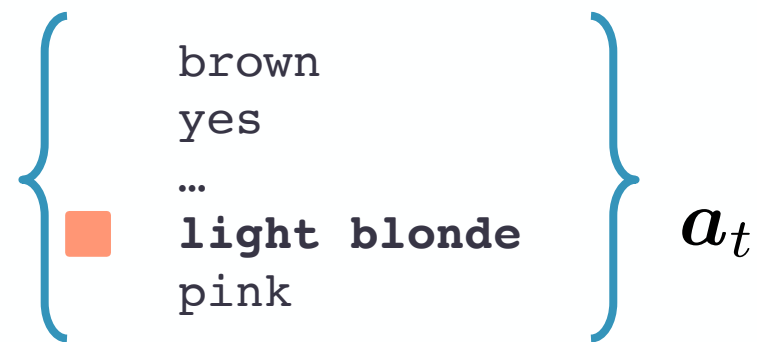
What color hair does the girl have?



# IVD EVALUATION

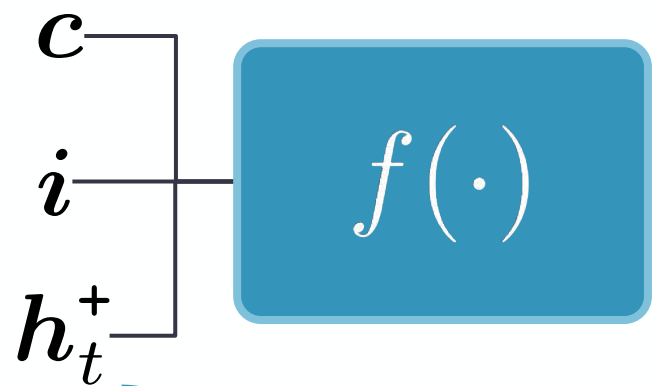


What color hair does the girl have?

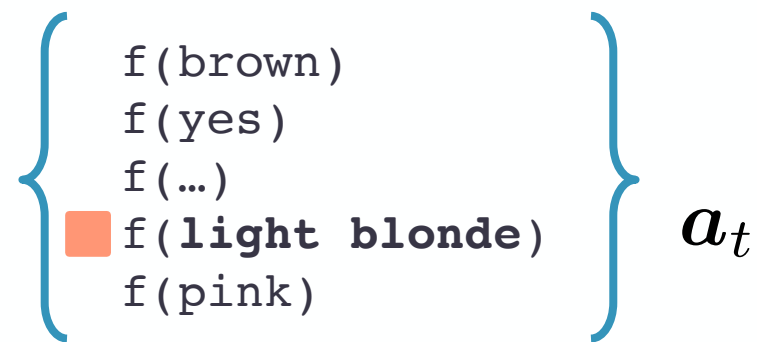




# IVD EVALUATION

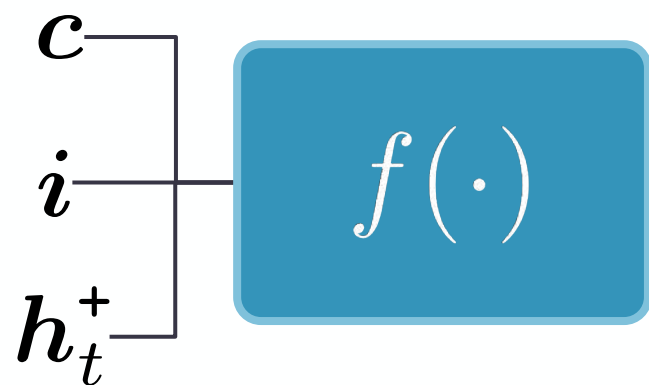


What color hair does the girl have?





# IVD EVALUATION

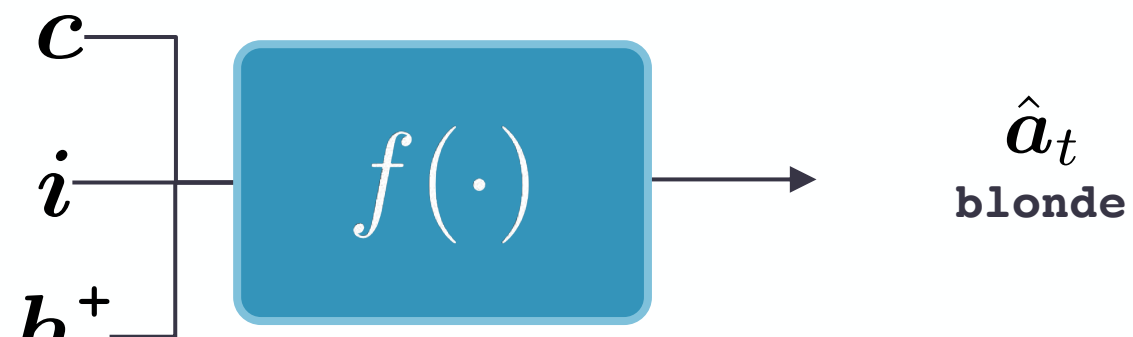


What color hair does the girl have?

$\left\{ \begin{array}{l} \text{f(light blonde)} \\ \text{f(light colour)} \\ \text{f(...)} \\ \text{f(none)} \\ \text{f(maybe)} \end{array} \right\} a_t$



# IVD EVALUATION



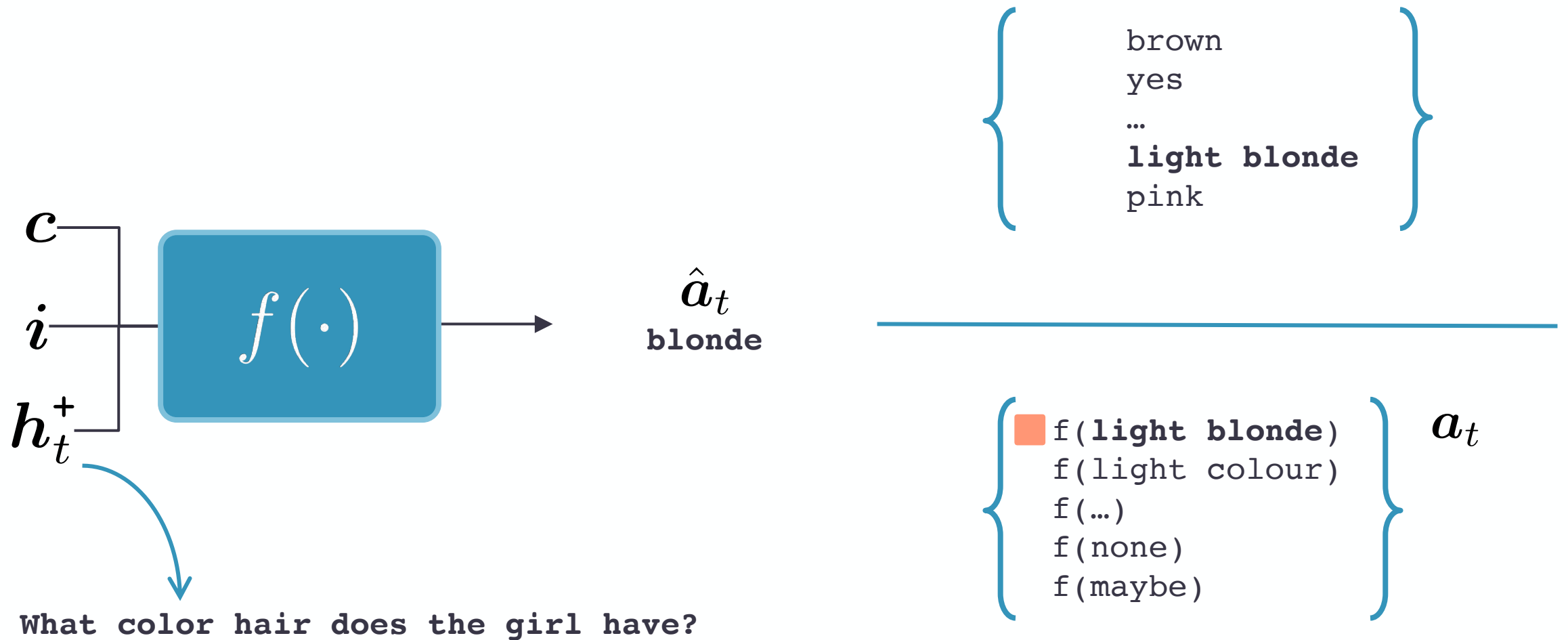
$h_t^+$

What color hair does the girl have?

$\left\{ \begin{array}{l} \text{f(light blonde)} \\ \text{f(light colour)} \\ \text{f(...)} \\ \text{f(none)} \\ \text{f(maybe)} \end{array} \right\} a_t$

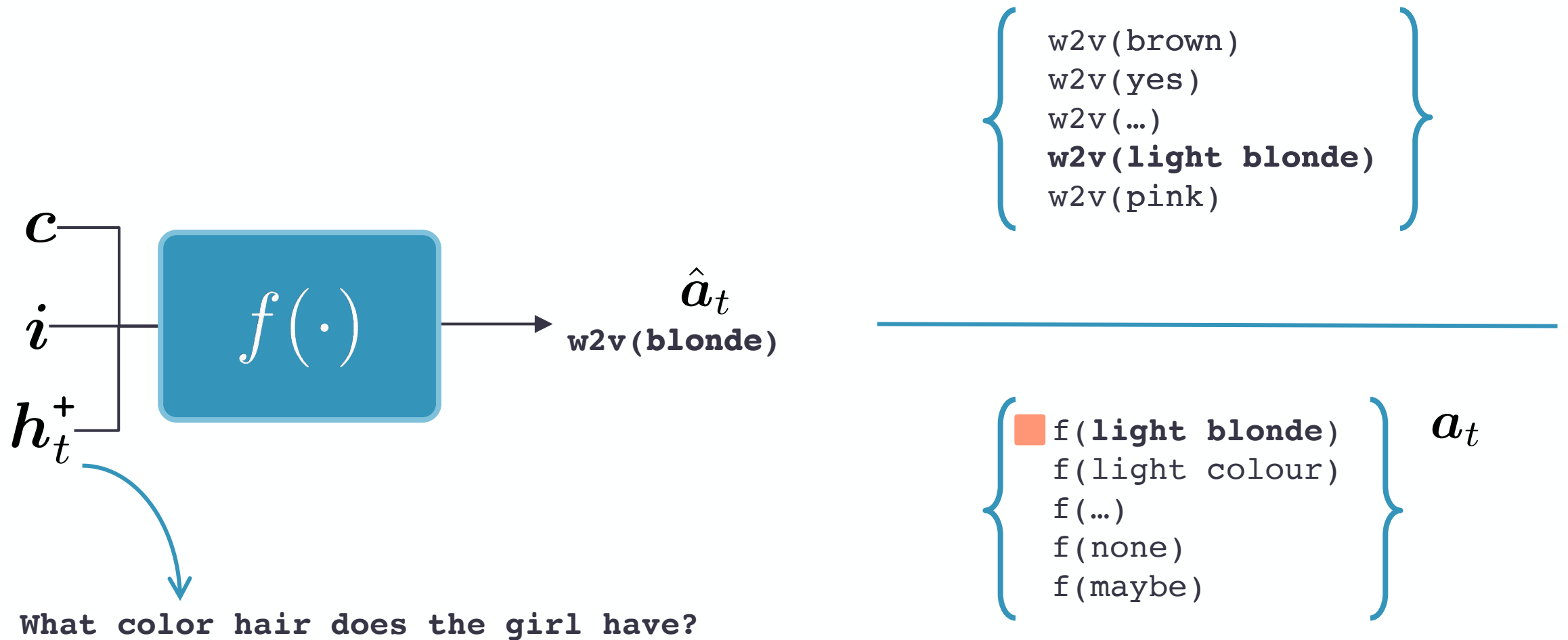


# IVD EVALUATION



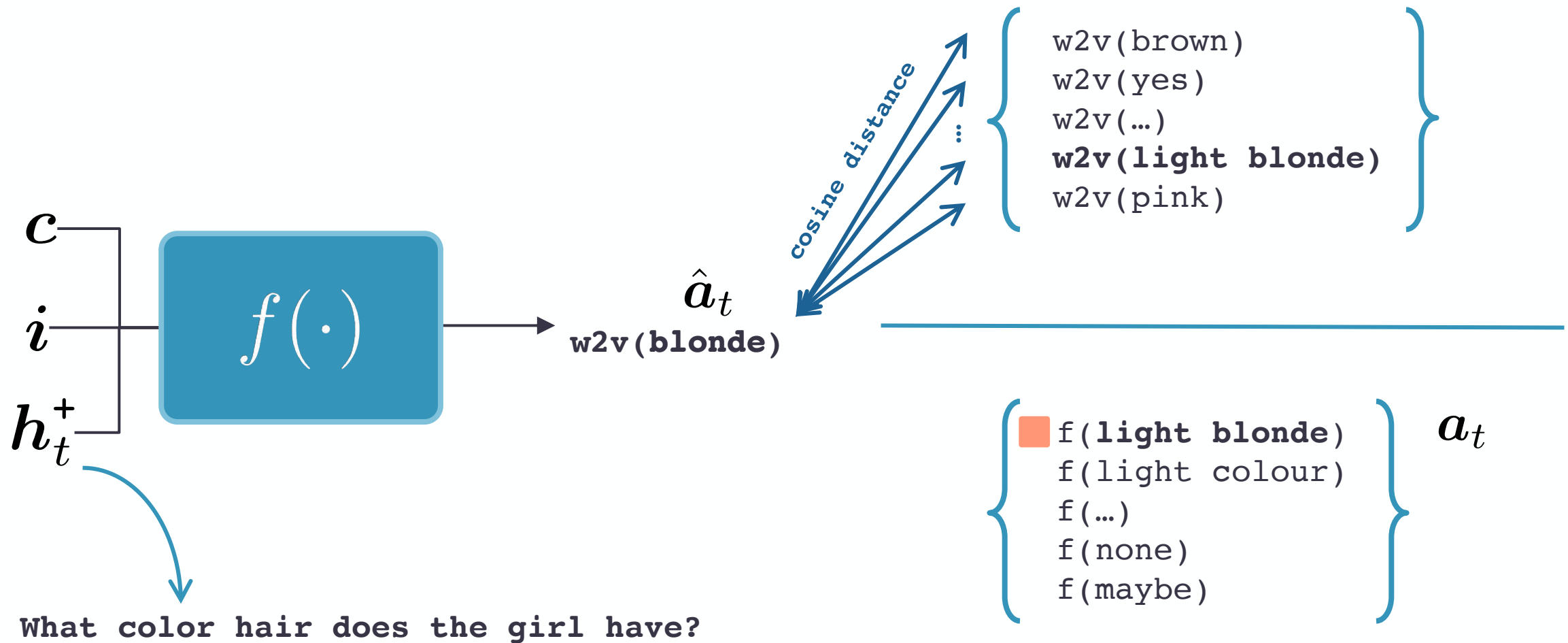


# IVD EVALUATION

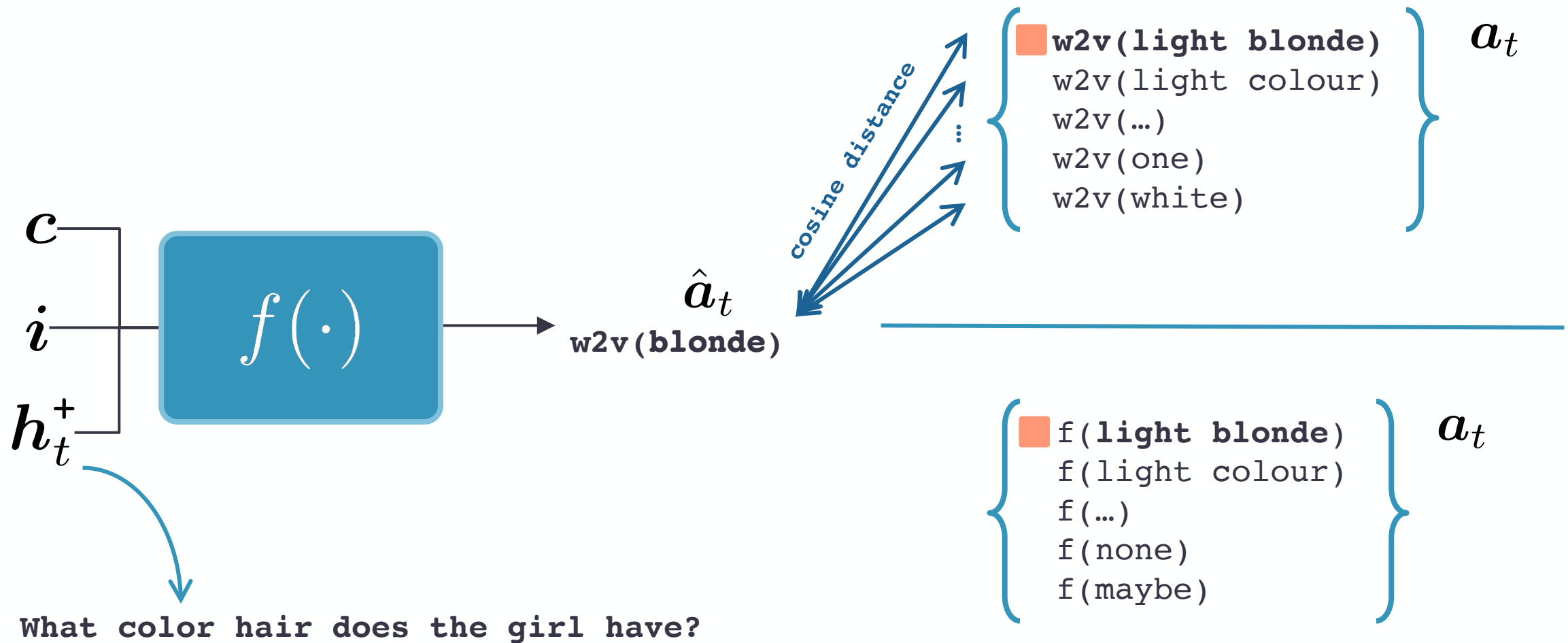




# IVD EVALUATION



# IVD EVALUATION





## IVD EVALUATION (I)

| Rank with...      |              | Mean Rank    |
|-------------------|--------------|--------------|
| <i>model</i>      | RL-QAbot [1] | 21.13        |
|                   | MN-QIH-G [2] | <b>17.06</b> |
|                   | <b>Ours</b>  | 20.38        |
|                   |              |              |
| <i>w2v cosine</i> | MN-QIH-G [2] | 31.31        |
|                   | <b>Ours</b>  | <b>25.65</b> |

[1] Das et al. Learning cooperative visual dialog agents with deep reinforcement learning. In ICCV, 2017

[2] Das et al. Visual Dialog. In CVPR, 2017

# IVD EVALUATION (II)

How old is the man?

What race is the man?

Is he wearing a shirt?

What color?

What gender is the baby?

What is the baby wearing?

What color is the remote?

Is the chair wood?

What color is the chair?



A man sitting in a chair holding a baby who is chewing on a remote

He looks thirty

White

Yes

Grey

I

Shirt

White

No

Brown

Twenties

White

Yes

Grey

Boy

T-shirt has pants

Silver

No

White has checkered

Thirtys

I can't see

No

White

Male

Shirt and shirt

It is white and black

No it's a chair

Gray

I say teenager

Caucasian

Yes

Black

I

Looks like white

White

No

Brown with texture



## ONE-WAY VISUAL DIALOGUE (IVD)

How many people  
are there?

it's crowded

Where is the  
pavement?

on the left

Where are the  
balloons?

in the air

A street with people and balloons

## TWO-WAY VISUAL DIALOGUE (2VD)

How many people  
are there?

it's crowded

Where is the  
pavement?

on the left

Where are the  
balloons?

in the air

A street with people and balloons

# CONDITIONAL VAE FOR IVD



image  $i$

A girl in a pink shirt eating a frosted donut

caption  $c$

|                                      |      |
|--------------------------------------|------|
| How old does the girl appear to be?  | Four |
| Does the donut have sprinkles?       | No   |
| Are there any other people with her? | No   |
| What color hair does the girl have?  |      |

dialogue history  $h_t^+ = \{h_{t-1}, q_t\}$

$$\log p_{\theta}(a_t \mid i, c, h_t^+)$$

Blonde



# CONDITIONAL VAE FOR 2VD




image  $i$

A girl in a pink shirt eating a frosted donut

caption  $c$

|                                      |        |
|--------------------------------------|--------|
| How old does the girl appear to be?  | Four   |
| Does the donut have sprinkles?       | No     |
| Are there any other people with her? | No     |
| What color hair does the girl have?  | Blonde |

dialogue  $d_{1:T}$

$$\log p_{\theta}(\textcolor{blue}{d}_{1:T} \mid \textcolor{brown}{i}, \textcolor{green}{c}, \textcolor{purple}{h})$$


# CONDITIONAL VAE FOR 2VD



image  $i$

A girl in a pink shirt eating a frosted donut

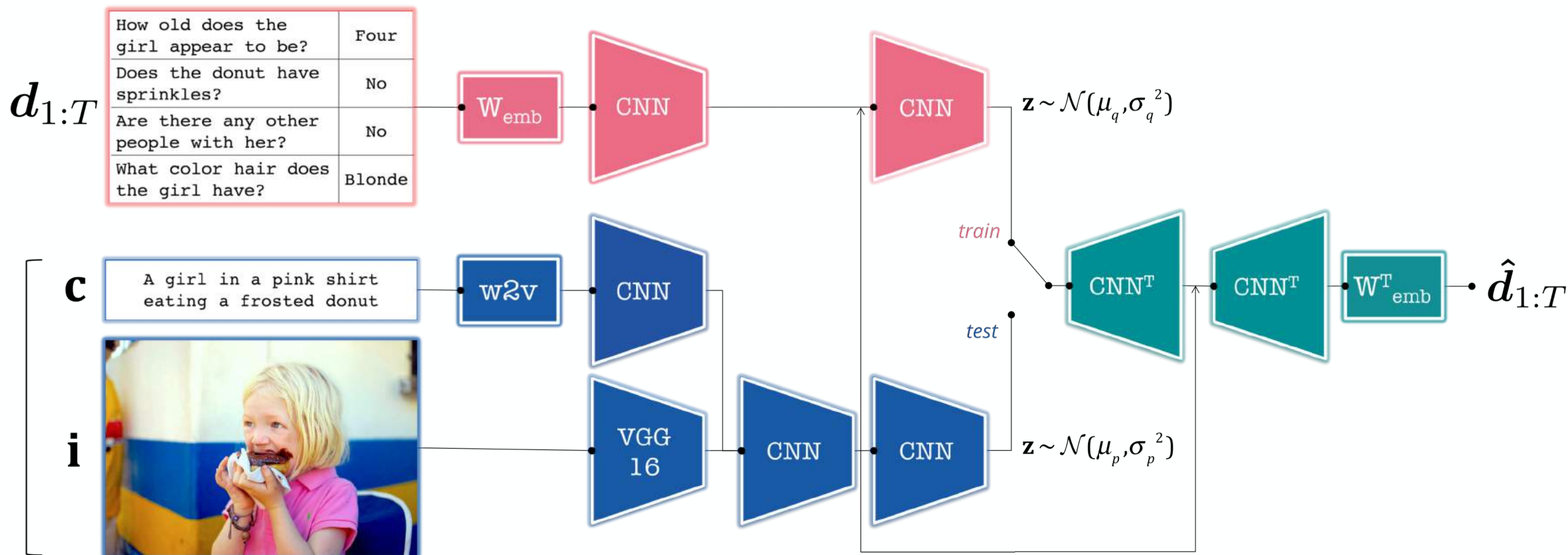
caption  $c$

|  |  |
|--|--|
|  |  |
|  |  |
|  |  |
|  |  |

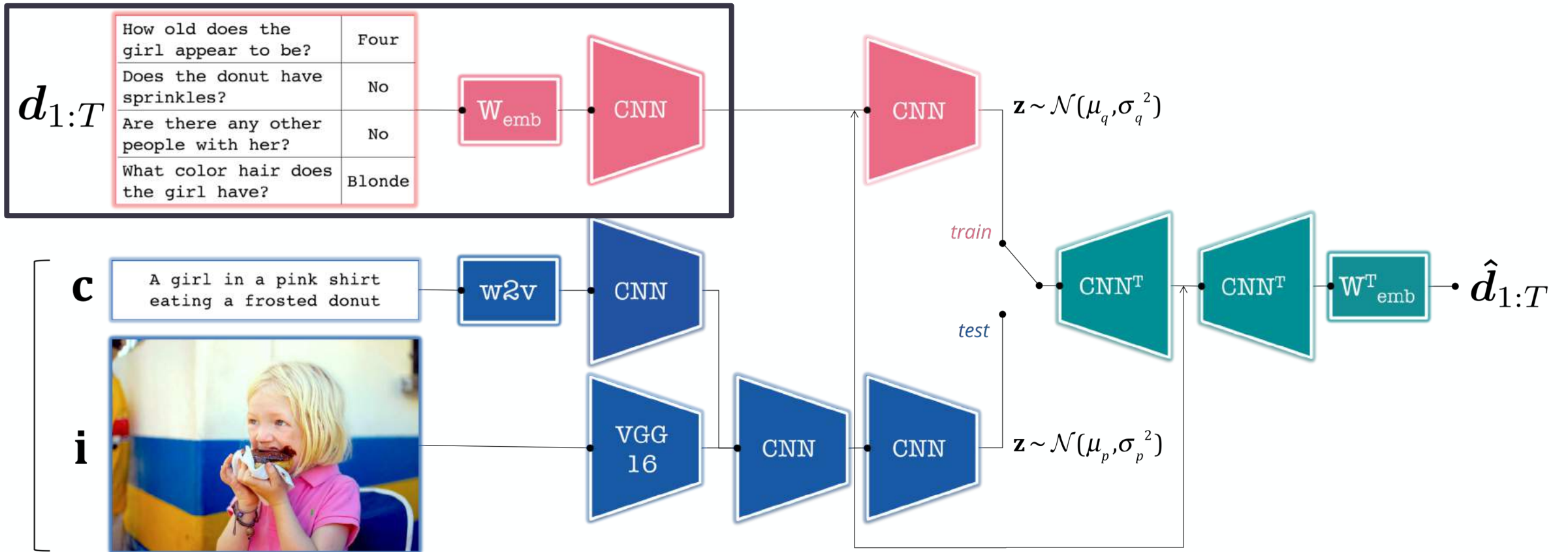
dialogue history  $h = \emptyset$

$$\log p_{\theta}(\underset{\uparrow}{d}_{1:T} \mid \underset{\uparrow}{i}, \underset{\uparrow}{c}, \underset{\uparrow}{h}) \geq \mathbb{E}_{q_{\phi}(z \mid \underset{\uparrow}{d}_{1:T}, \underset{\uparrow}{i}, \underset{\uparrow}{c}, \underset{\uparrow}{h})} [\log p_{\theta}(\underset{\uparrow}{d}_{1:T} \mid z, \underset{\uparrow}{i}, \underset{\uparrow}{c}, \underset{\uparrow}{h})] - \mathbb{D}_{\text{KL}}(q_{\phi}(z \mid \underset{\uparrow}{d}_{1:T}, \underset{\uparrow}{i}, \underset{\uparrow}{c}, \underset{\uparrow}{h}) \parallel p_{\theta}(z \mid \underset{\uparrow}{i}, \underset{\uparrow}{c}, \underset{\uparrow}{h}))$$





$$\log p_{\theta}(\mathbf{d}_{1:T} \mid \mathbf{i}, \mathbf{c}, \mathbf{h}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{d}_{1:T}, \mathbf{i}, \mathbf{c}, \mathbf{h})} [\log p_{\theta}(\mathbf{d}_{1:T} \mid \mathbf{z}, \mathbf{i}, \mathbf{c}, \mathbf{h})] - \mathbb{D}_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{d}_{1:T}, \mathbf{i}, \mathbf{c}, \mathbf{h}) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{i}, \mathbf{c}, \mathbf{h}))$$



$$\log p_{\theta}(d_{1:T} \mid i, c, h) \geq \mathbb{E}_{q_{\phi}(z \mid d_{1:T}, i, c, h)} [\log p_{\theta}(d_{1:T} \mid z, i, c, h)] - \mathbb{D}_{\text{KL}}(q_{\phi}(z \mid d_{1:T}, i, c, h) \parallel p_{\theta}(z \mid i, c, h))$$



# ‘COLOURING’ DIALOGUE WITH CONVOLUTIONS

|                                      |        |
|--------------------------------------|--------|
| How old does the girl appear to be?  | Four   |
| Does the donut have sprinkles?       | No     |
| Are there any other people with her? | No     |
| What color hair does the girl have?  | Blonde |

# 'COLOURING' DIALOGUE WITH CONVOLUTIONS

How old does the girl appear to be?

Four

Does the donut have sprinkles?

No

Are there any other people with her?

No

What color hair does the girl have?

Blonde



# 'COLOURING' DIALOGUE WITH CONVOLUTIONS



How old does the girl appear to be?

Four

Does the donut have sprinkles?

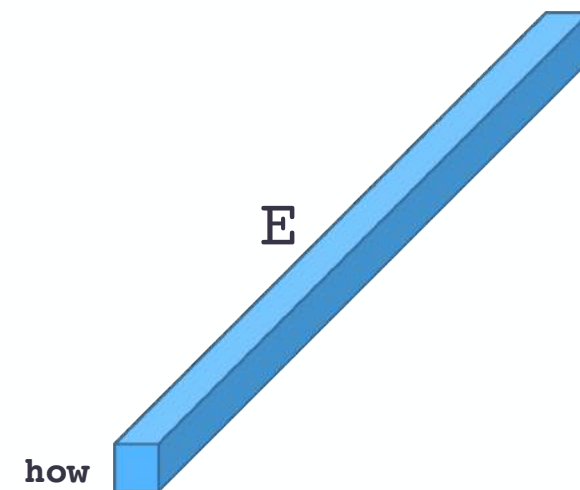
No

Are there any other people with her?

No

What color hair does the girl have?

Blonde



# 'COLOURING' DIALOGUE WITH CONVOLUTIONS



How old does the girl appear to be?

Four

Does the donut have sprinkles?

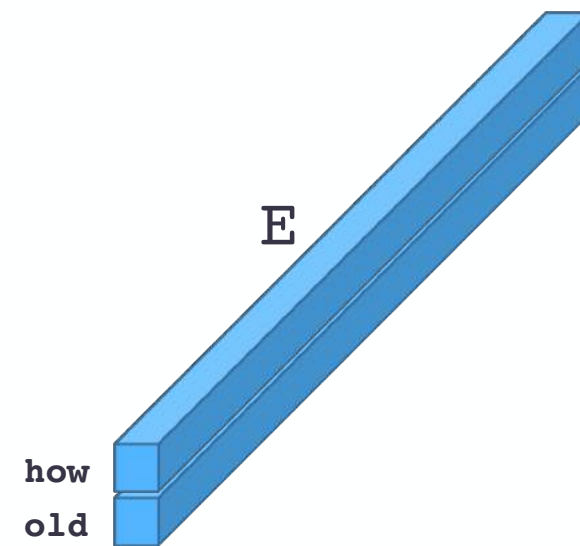
No

Are there any other people with her?

No

What color hair does the girl have?

Blonde





# 'COLOURING' DIALOGUE WITH CONVOLUTIONS



How old does the girl appear to be?

Four

Does the donut have sprinkles?

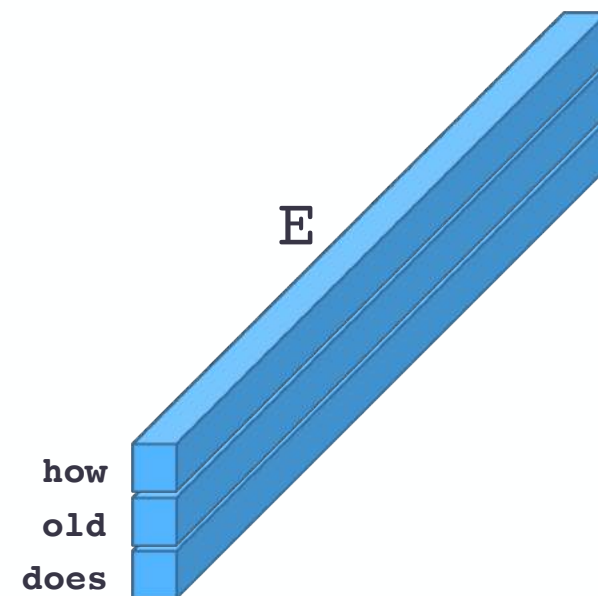
No

Are there any other people with her?

No

What color hair does the girl have?

Blonde



# 'COLOURING' DIALOGUE WITH CONVOLUTIONS



How old does the girl appear to be?

Four

Does the donut have sprinkles?

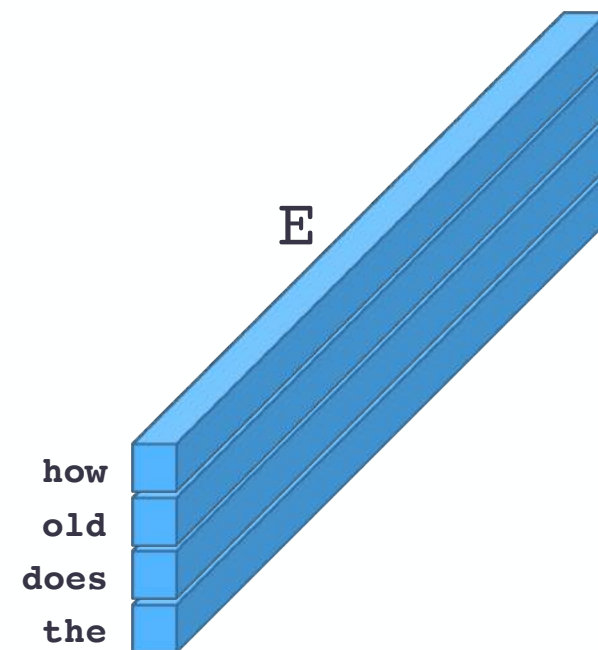
No

Are there any other people with her?

No

What color hair does the girl have?

Blonde





# 'COLOURING' DIALOGUE WITH CONVOLUTIONS



How old does the girl appear to be?

Four

Does the donut have sprinkles?

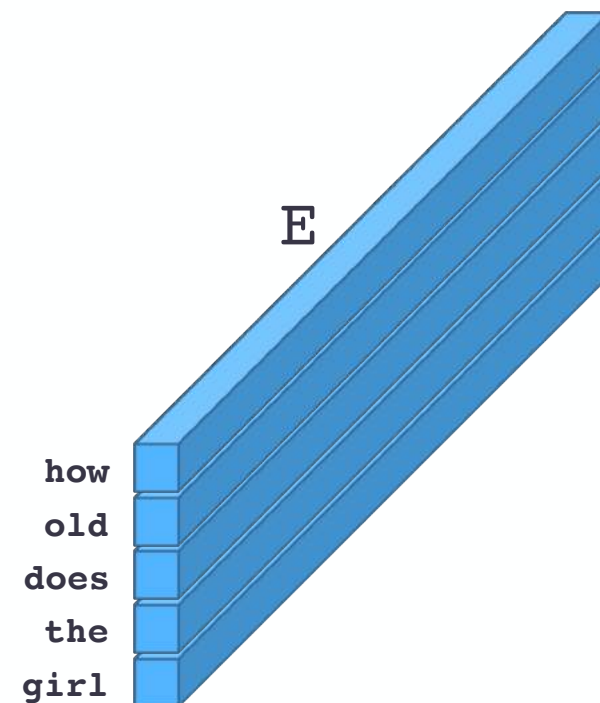
No

Are there any other people with her?

No

What color hair does the girl have?

Blonde



# 'COLOURING' DIALOGUE WITH CONVOLUTIONS



How old does the girl **appear** to be?

Four

Does the donut have sprinkles?

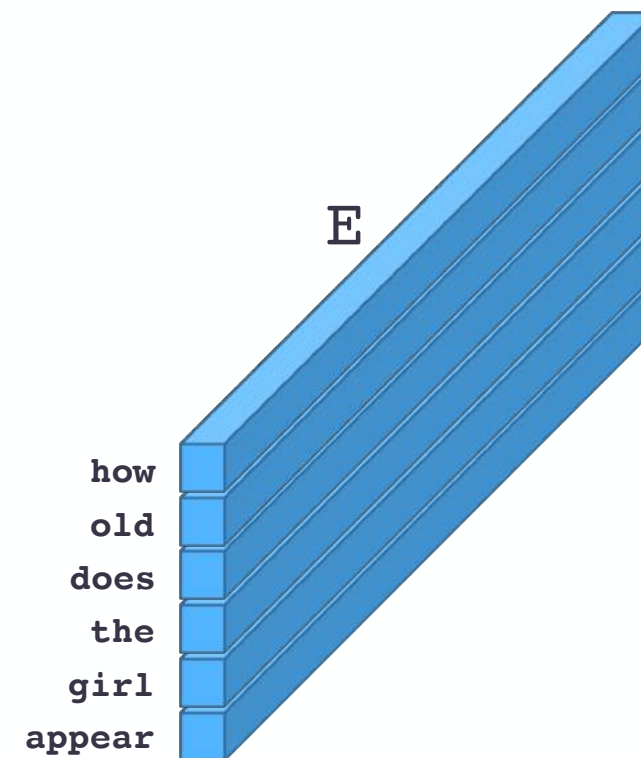
No

Are there any other people with her?

No

What color hair does the girl have?

Blonde





# 'COLOURING' DIALOGUE WITH CONVOLUTIONS



How old does the girl appear **to** be?

Four

Does the donut have sprinkles?

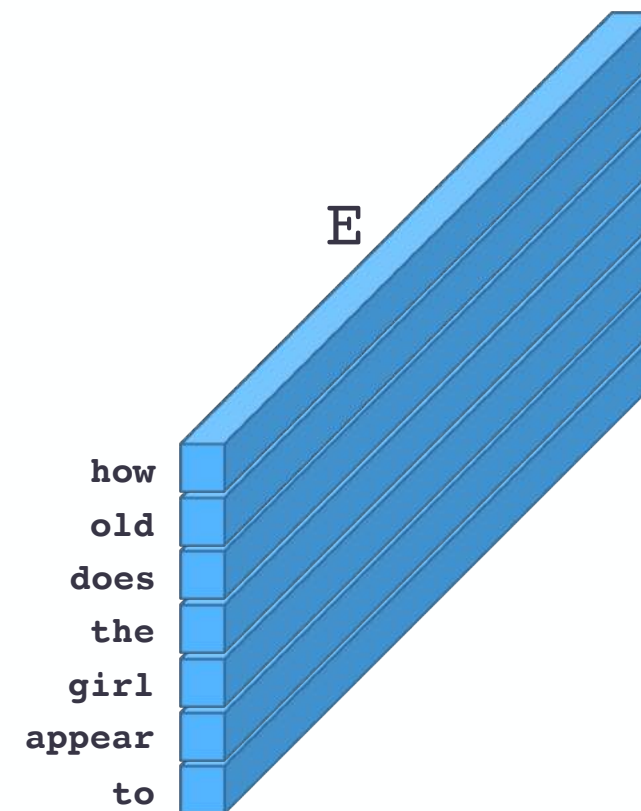
No

Are there any other people with her?

No

What color hair does the girl have?

Blonde



# 'COLOURING' DIALOGUE WITH CONVOLUTIONS



How old does the girl appear to be?

Four

Does the donut have sprinkles?

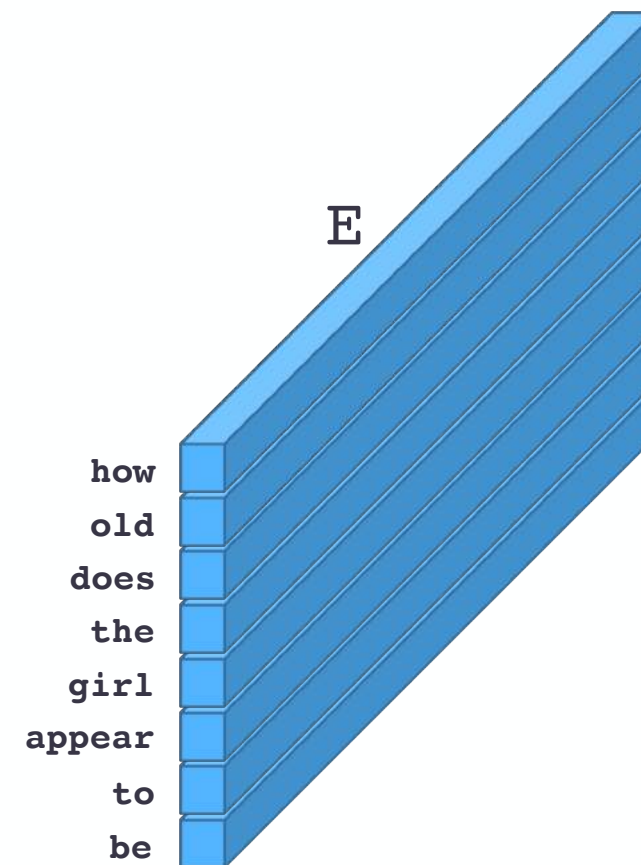
No

Are there any other people with her?

No

What color hair does the girl have?

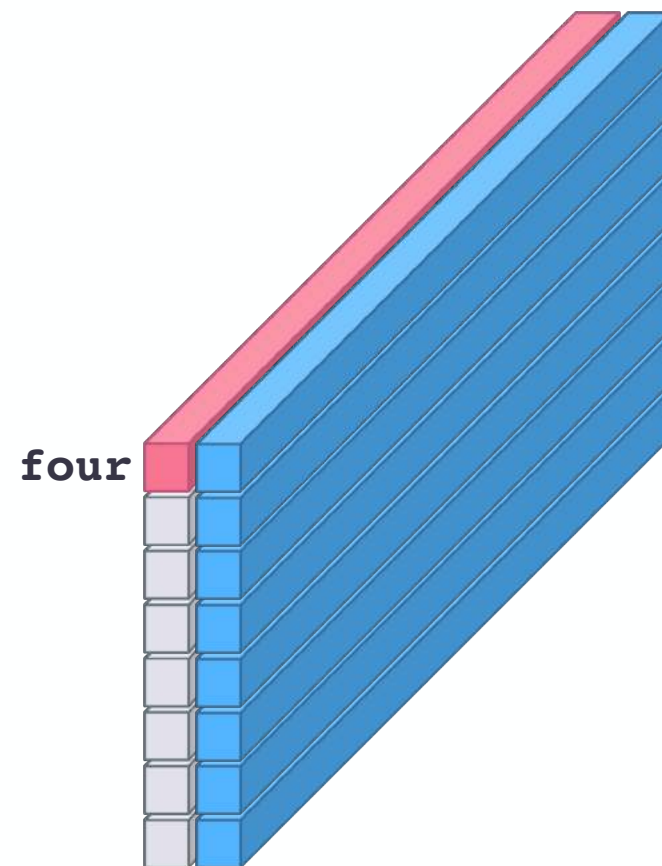
Blonde





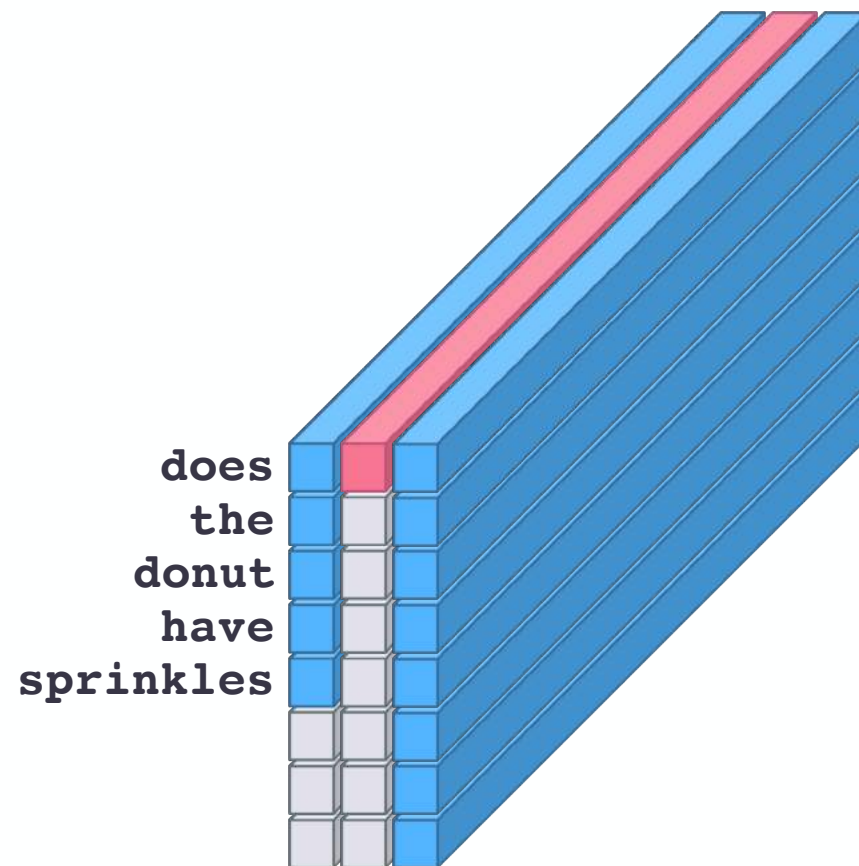
# 'COLOURING' DIALOGUE WITH CONVOLUTIONS

|   |                                      |
|---|--------------------------------------|
|   | How old does the girl appear to be?  |
| → | Four                                 |
|   | Does the donut have sprinkles?       |
|   | No                                   |
|   | Are there any other people with her? |
|   | No                                   |
|   | What color hair does the girl have?  |
|   | Blonde                               |



# 'COLOURING' DIALOGUE WITH CONVOLUTIONS

|   |                                      |
|---|--------------------------------------|
|   | How old does the girl appear to be?  |
|   | Four                                 |
| → | Does the donut have sprinkles?       |
|   | No                                   |
|   | Are there any other people with her? |
|   | No                                   |
|   | What color hair does the girl have?  |
|   | Blonde                               |





# 'COLOURING' DIALOGUE WITH CONVOLUTIONS

How old does the girl appear to be?

Four

Does the donut have sprinkles?

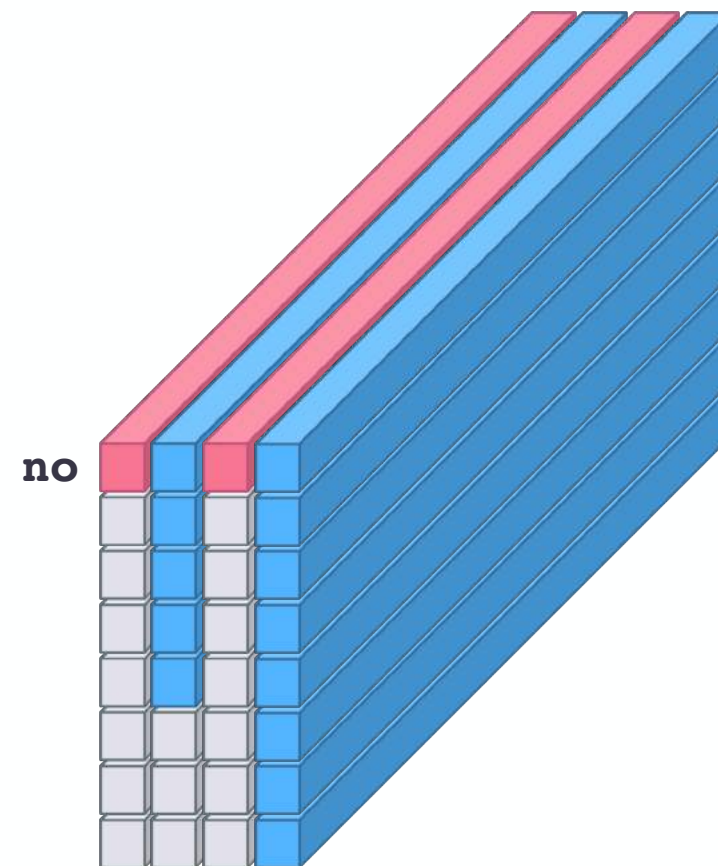
→ No

Are there any other people with her?

No

What color hair does the girl have?

Blonde



# 'COLOURING' DIALOGUE WITH CONVOLUTIONS

How old does the girl appear to be?

Four

Does the donut have sprinkles?

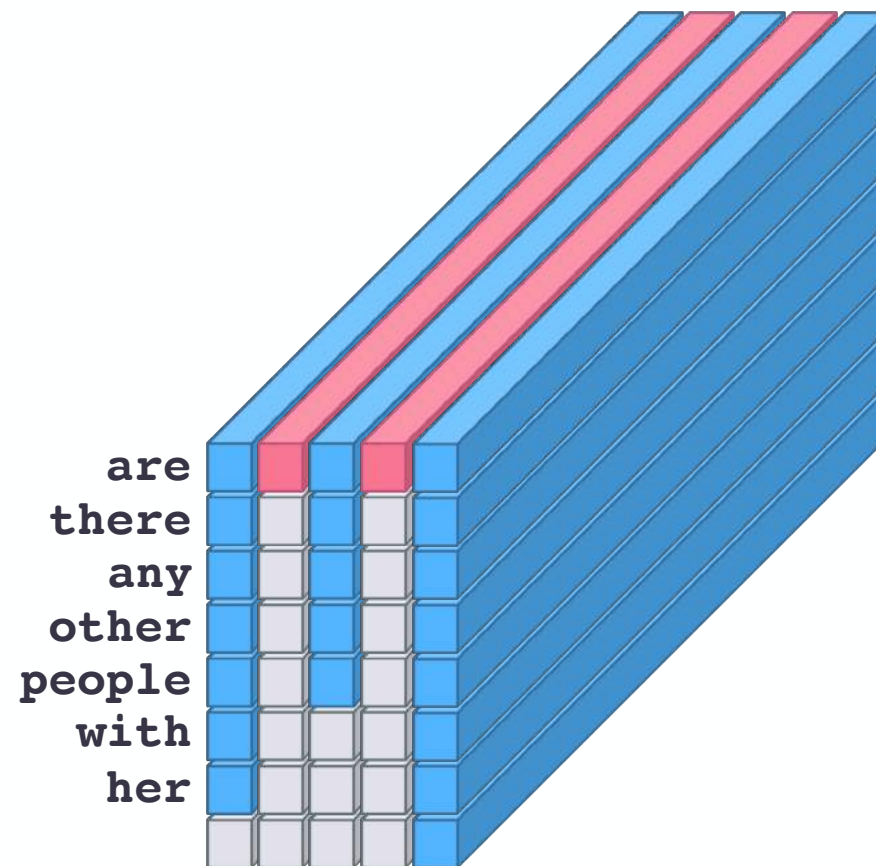
No

→ Are there any other people with her?

No

What color hair does the girl have?

Blonde





# 'COLOURING' DIALOGUE WITH CONVOLUTIONS

How old does the girl appear to be?

Four

Does the donut have sprinkles?

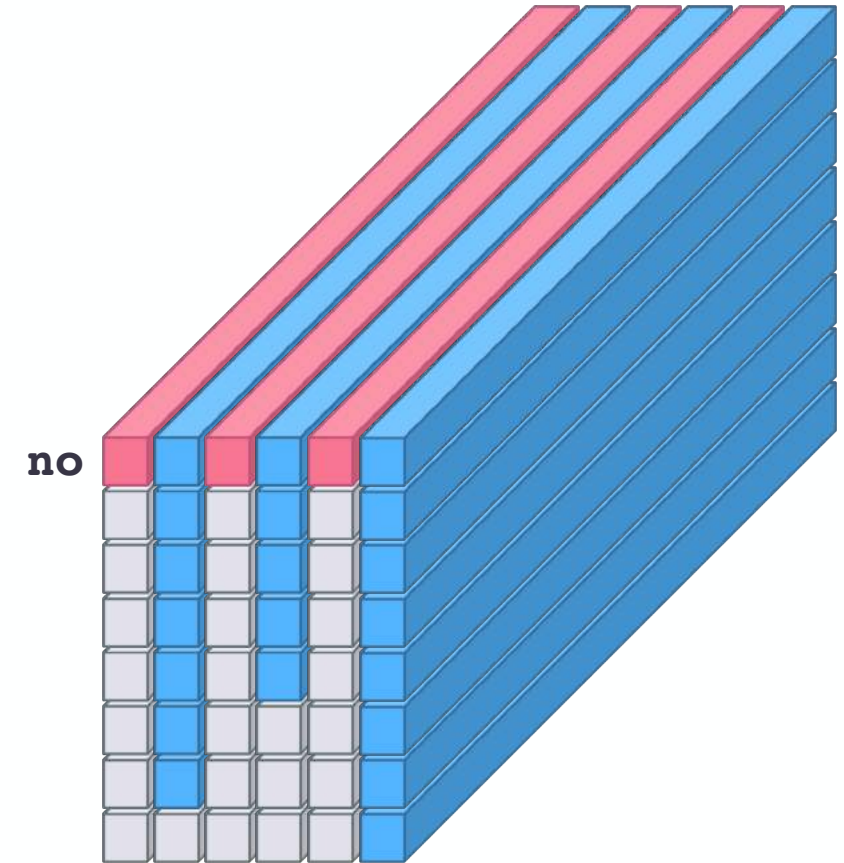
No

Are there any other people with her?

→ No

What color hair does the girl have?

Blonde



# 'COLOURING' DIALOGUE WITH CONVOLUTIONS

How old does the girl appear to be?

Four

Does the donut have sprinkles?

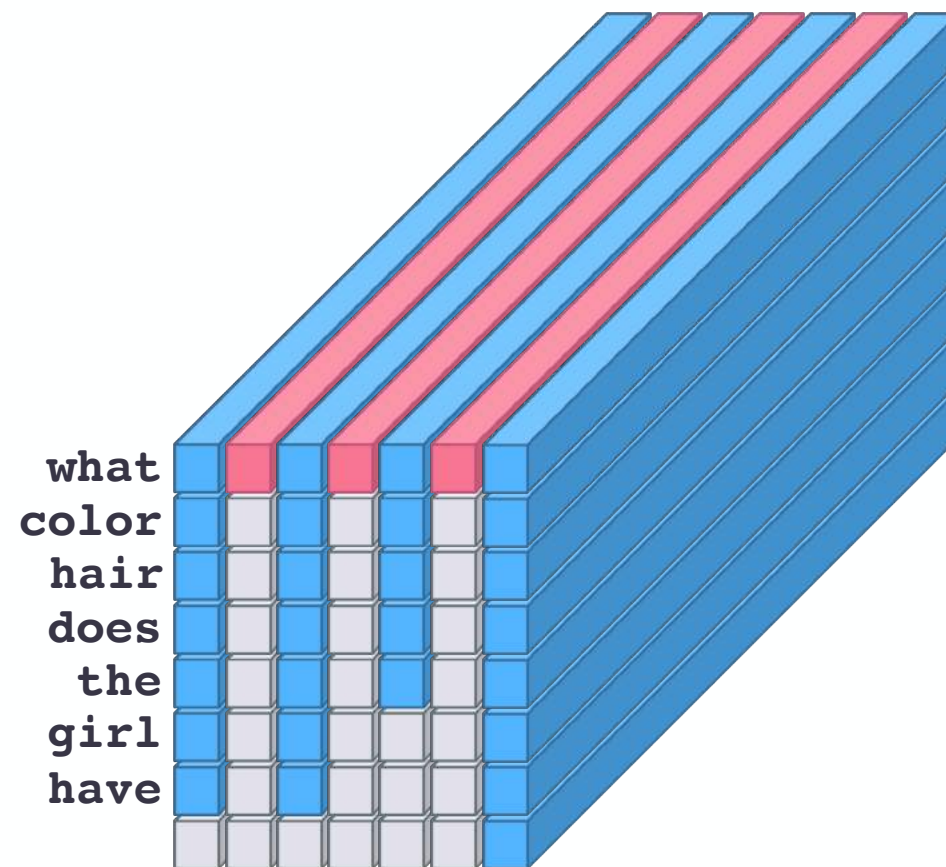
No

Are there any other people with her?

No

→ What color hair does the girl have?

Blonde



# 'COLOURING' DIALOGUE WITH CONVOLUTIONS

How old does the girl appear to be?

Four

Does the donut have sprinkles?

No

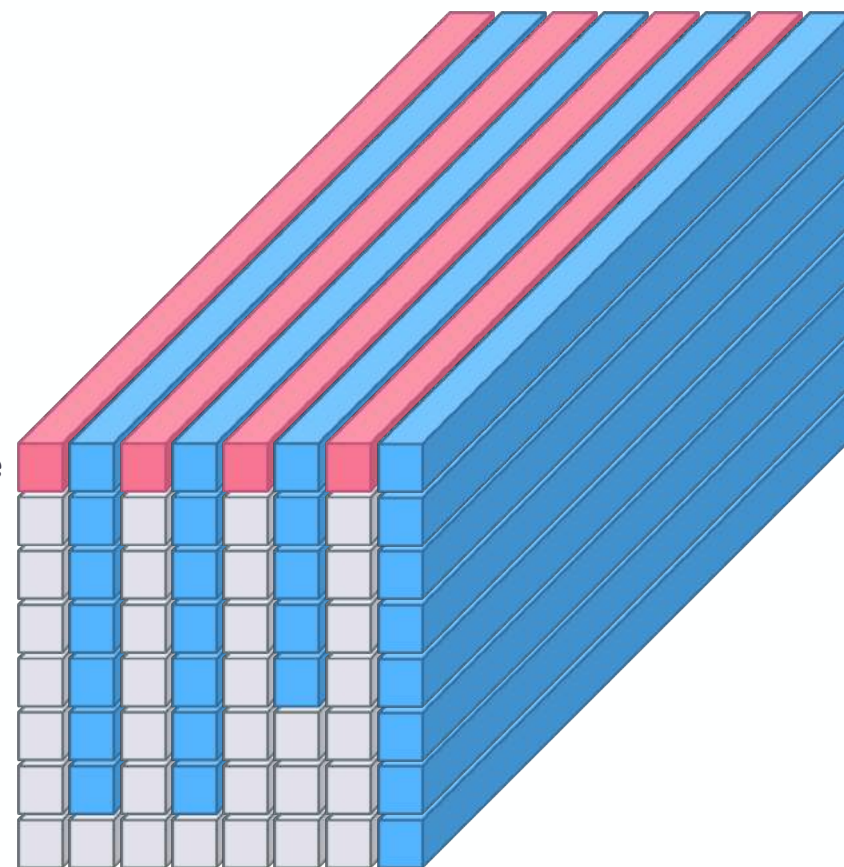
Are there any other people with her?

No

What color hair does the girl have?

→ Blonde

**blonde**





# 'COLOURING' DIALOGUE WITH CONVOLUTIONS

How old does the girl appear to be?

Four

Does the donut have sprinkles?

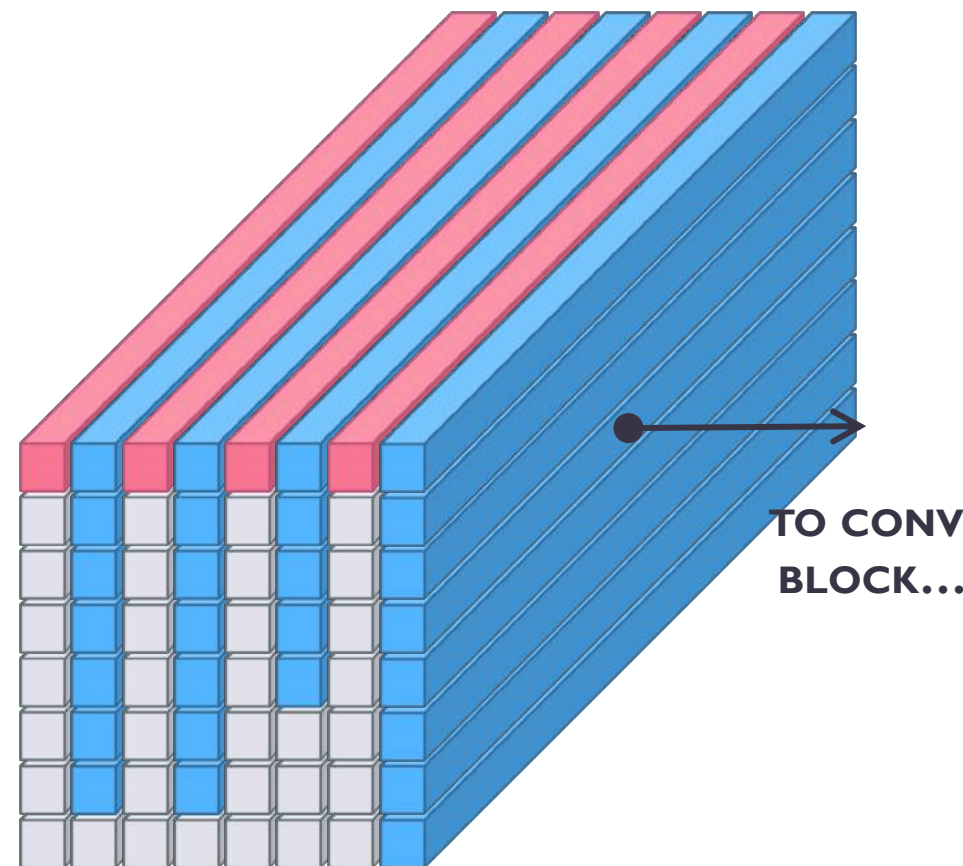
No

Are there any other people with her?

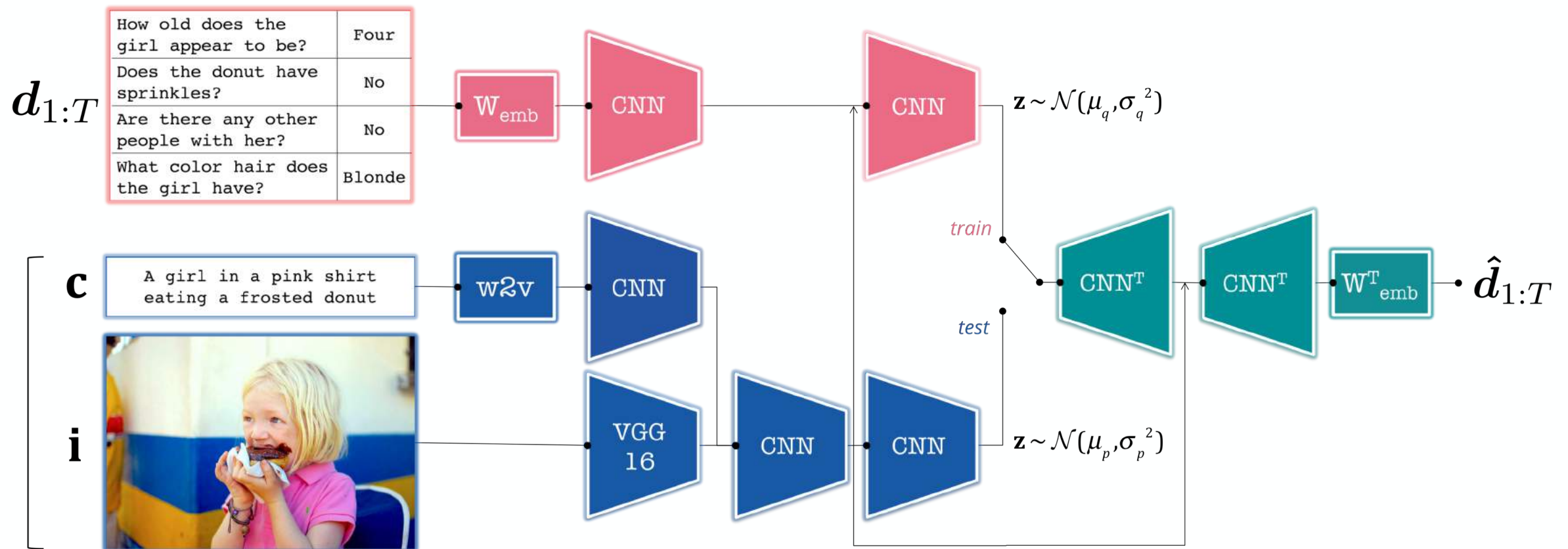
No

What color hair does the girl have?

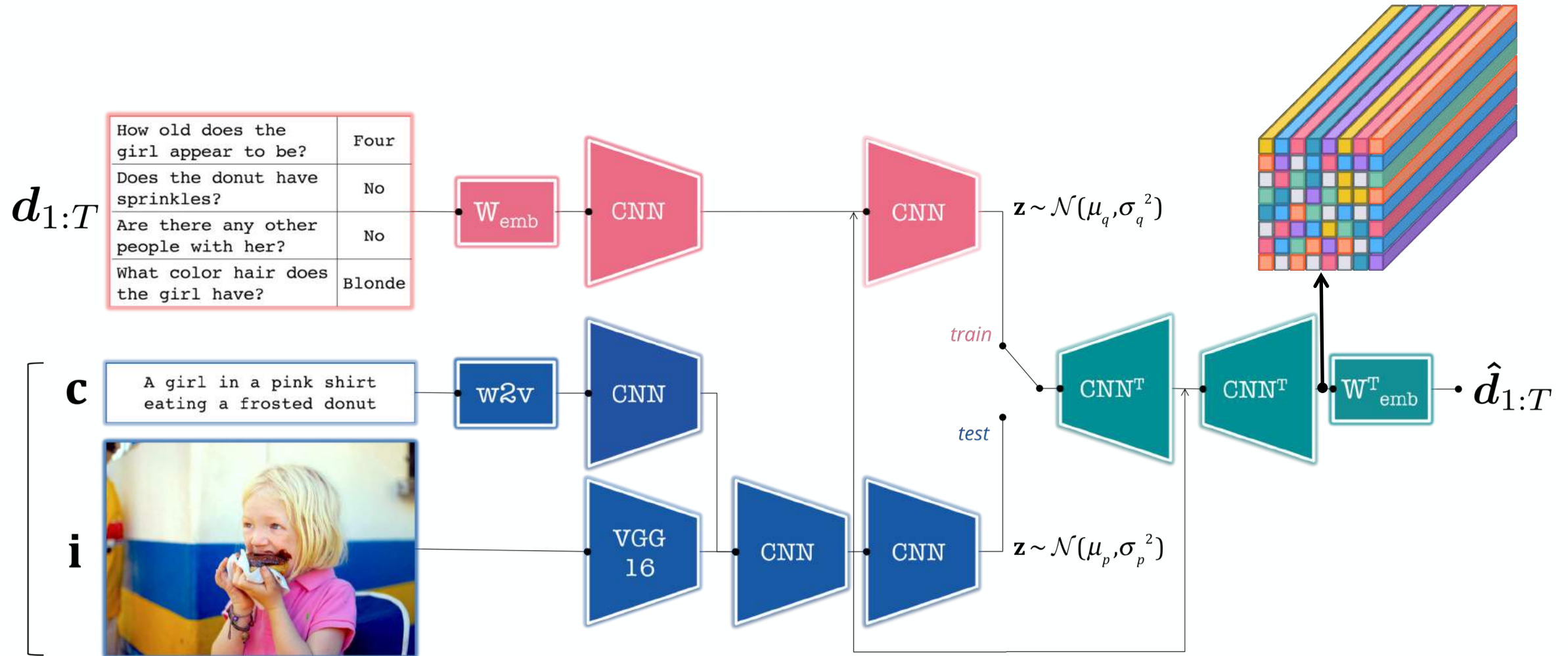
Blonde



# EXPLICIT SEQUENCE WITH AN AUTOREGRESSIVE DECODER

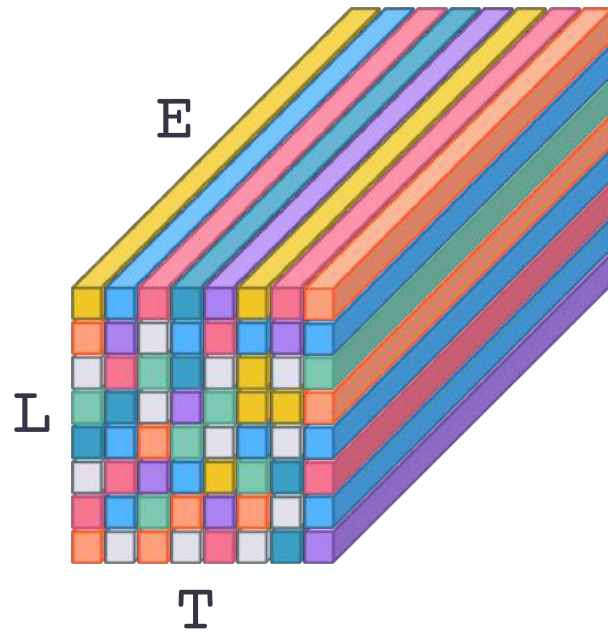


# EXPLICIT SEQUENCE WITH AN AUTOREGRESSIVE DECODER

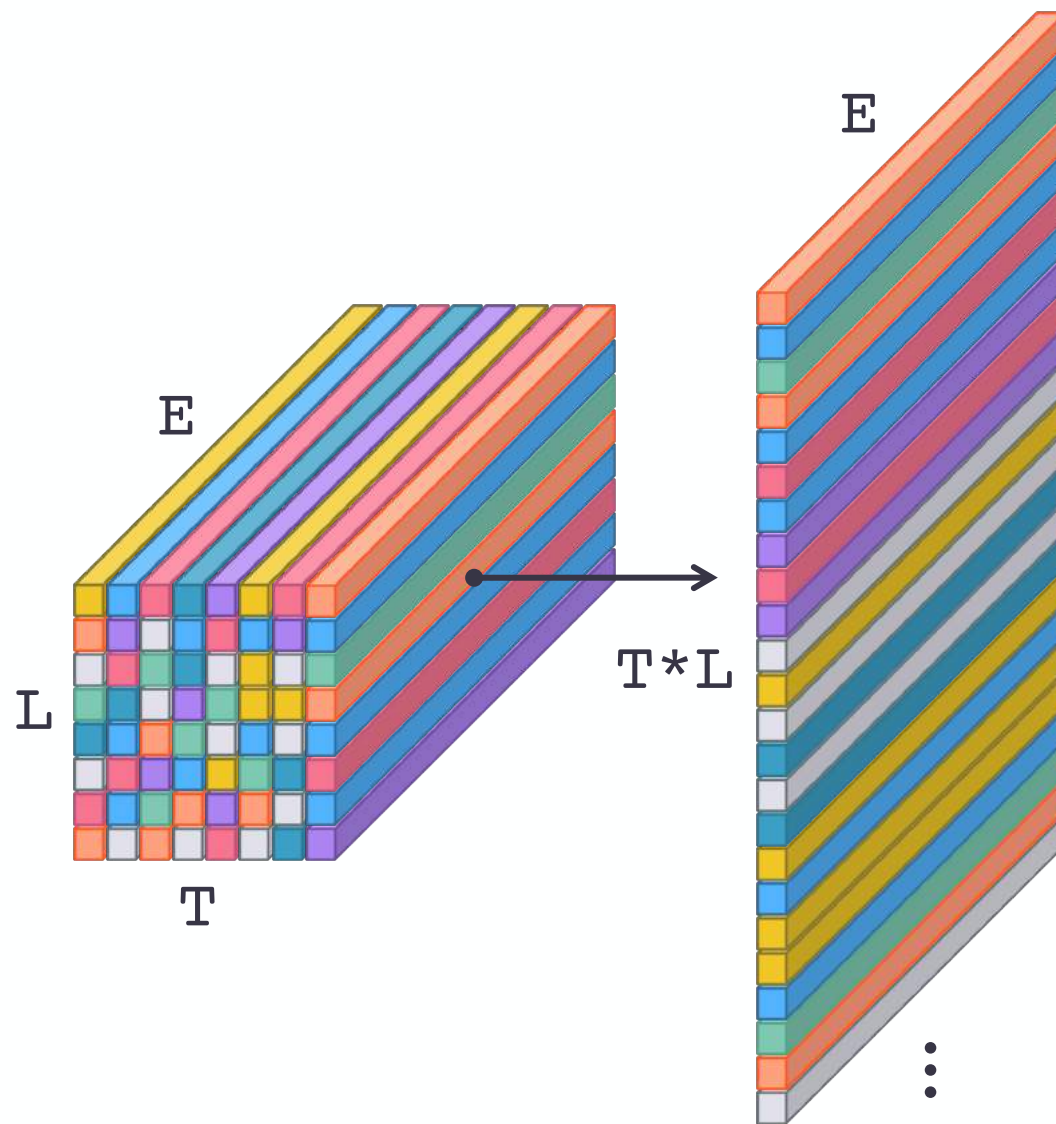




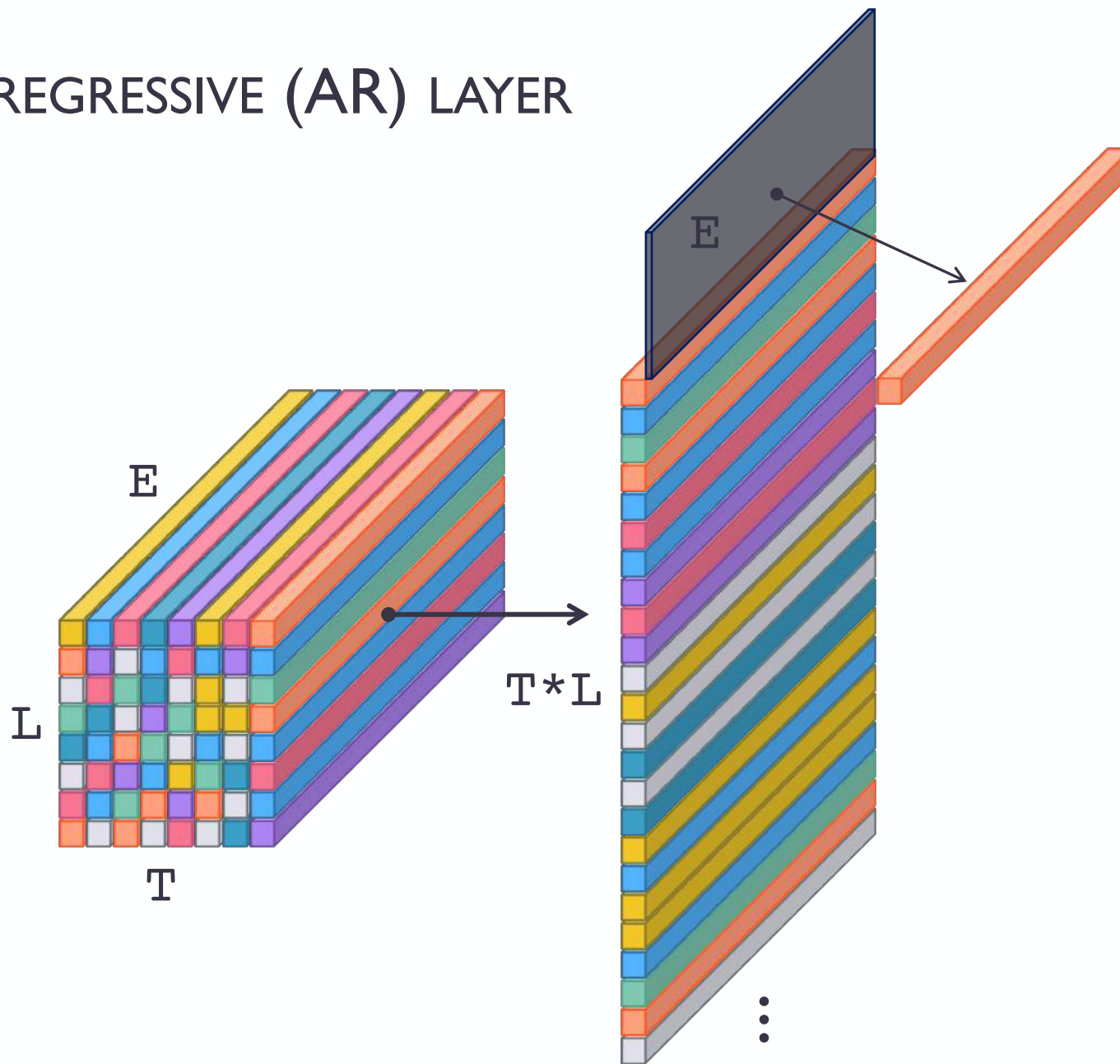
# AN AUTO-REGRESSIVE (AR) LAYER



# AN AUTO-REGRESSIVE (AR) LAYER

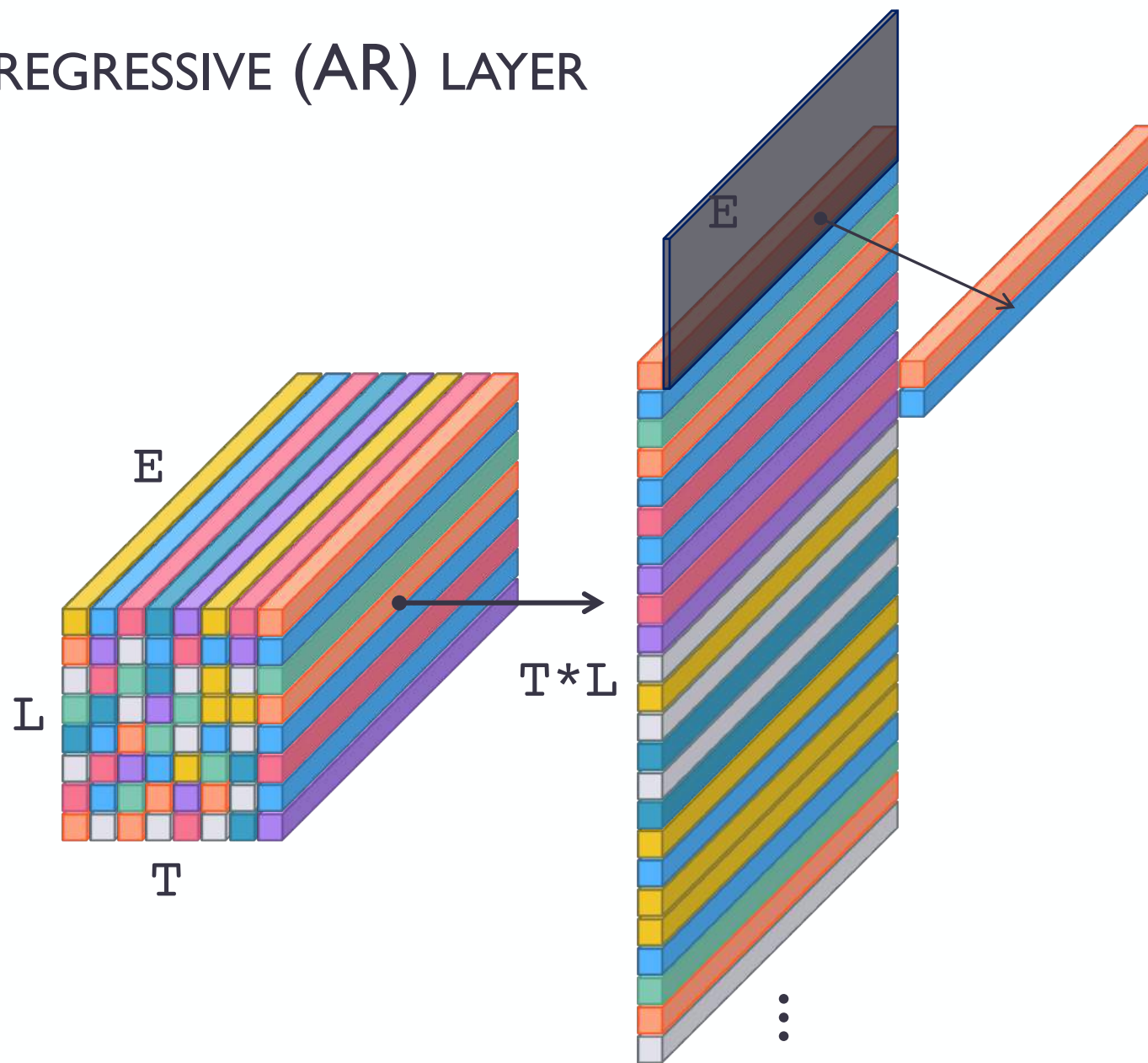


# AN AUTO-REGRESSIVE (AR) LAYER

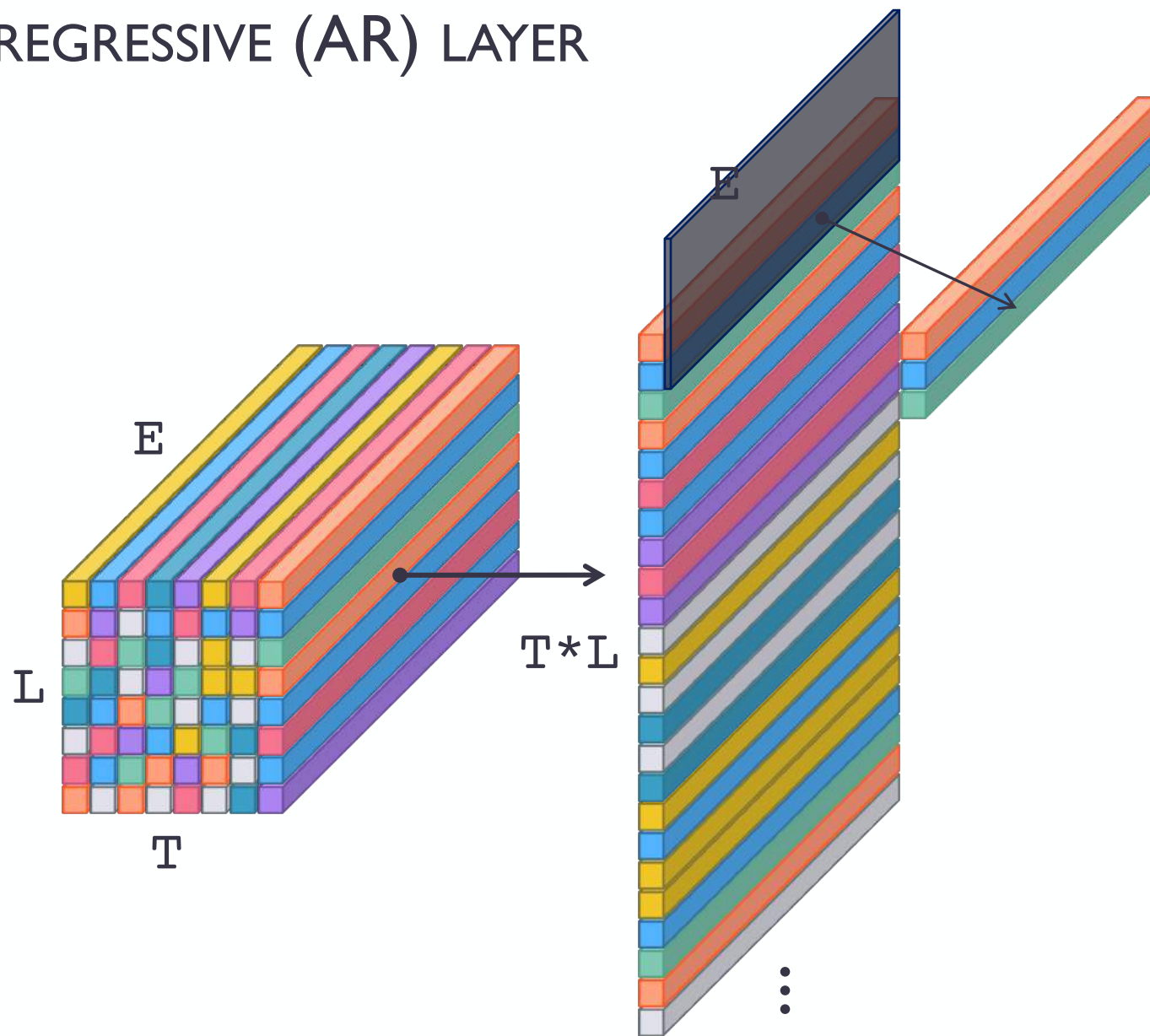




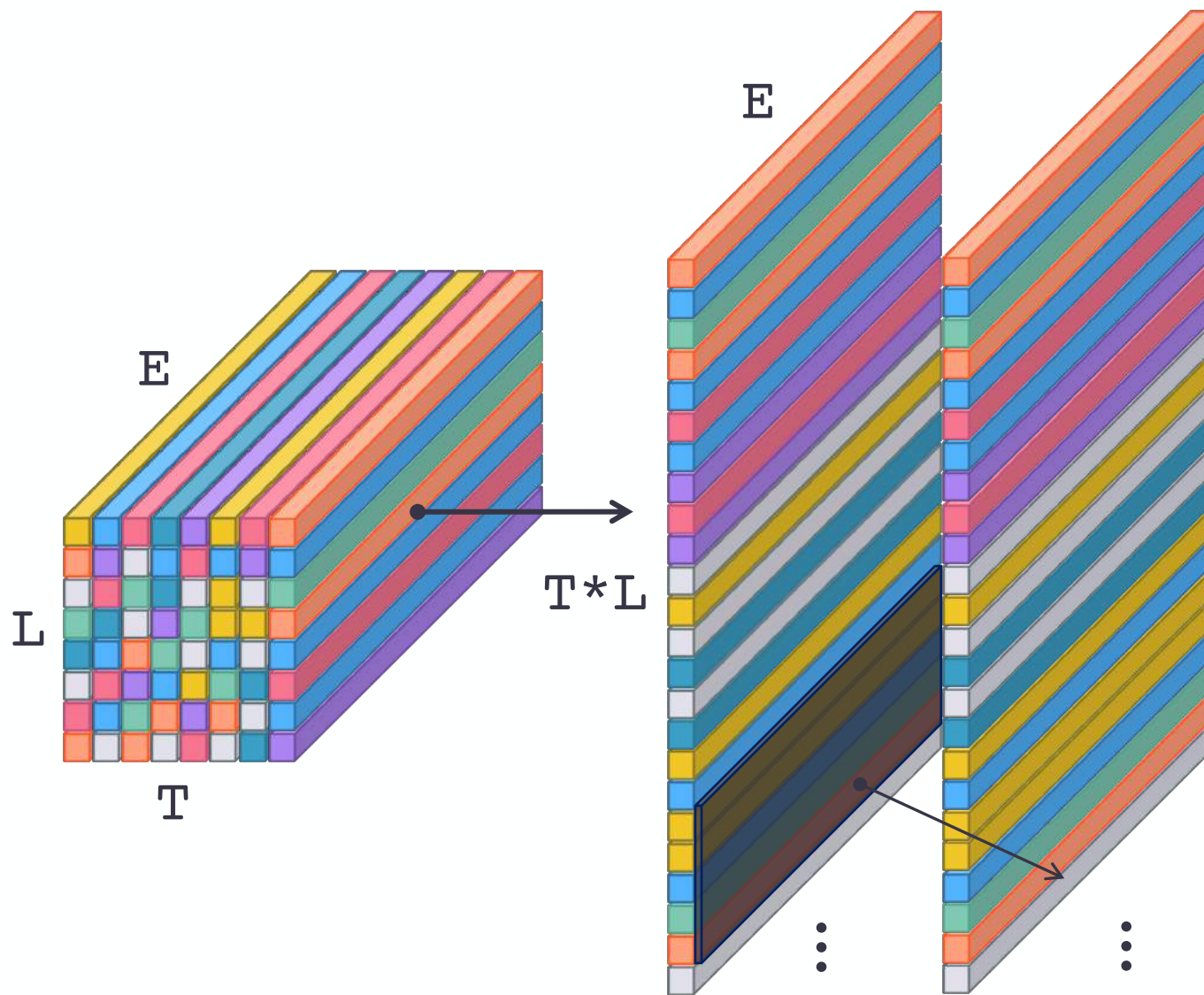
# AN AUTO-REGRESSIVE (AR) LAYER



# AN AUTO-REGRESSIVE (AR) LAYER

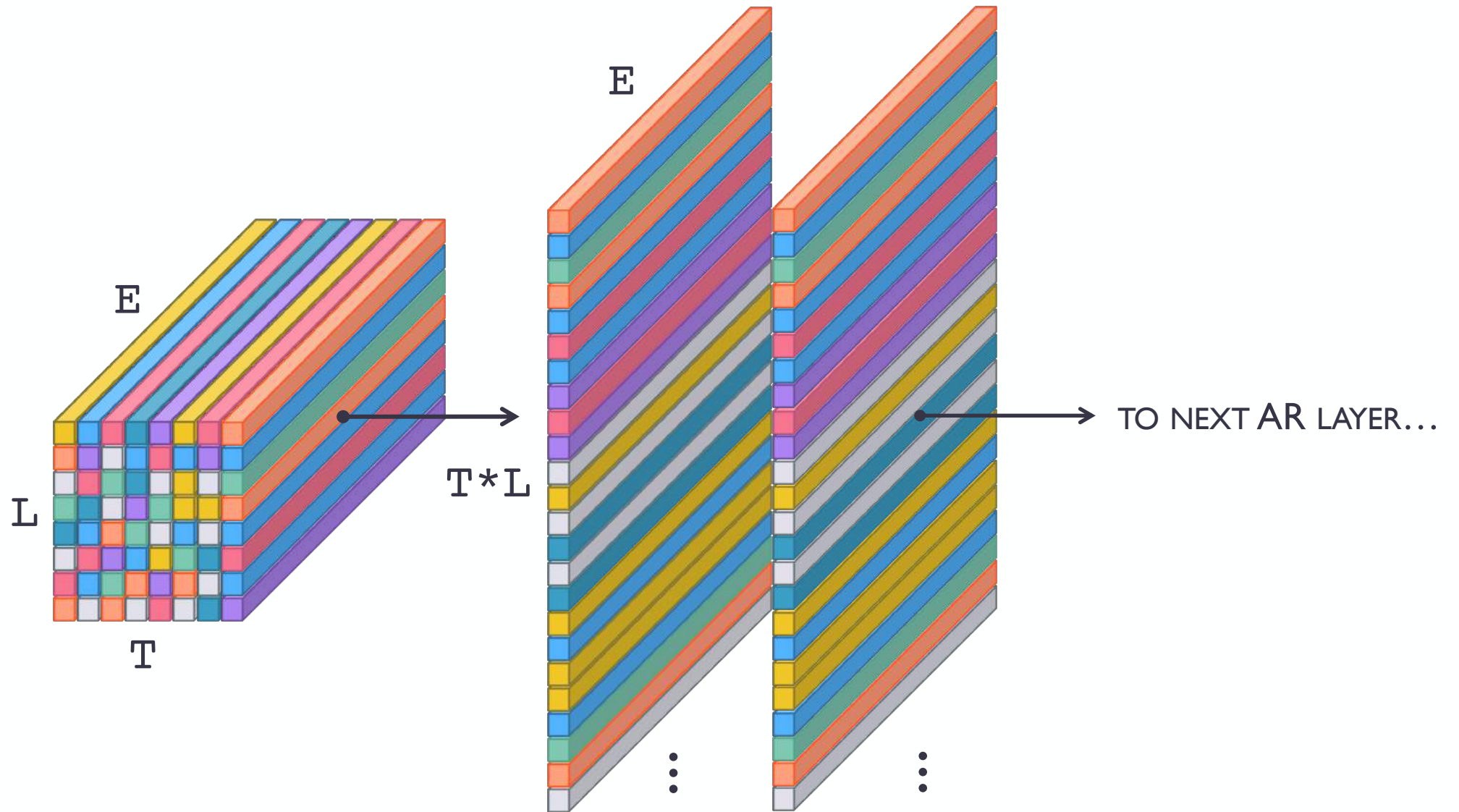


# AN AUTO-REGRESSIVE (AR) LAYER

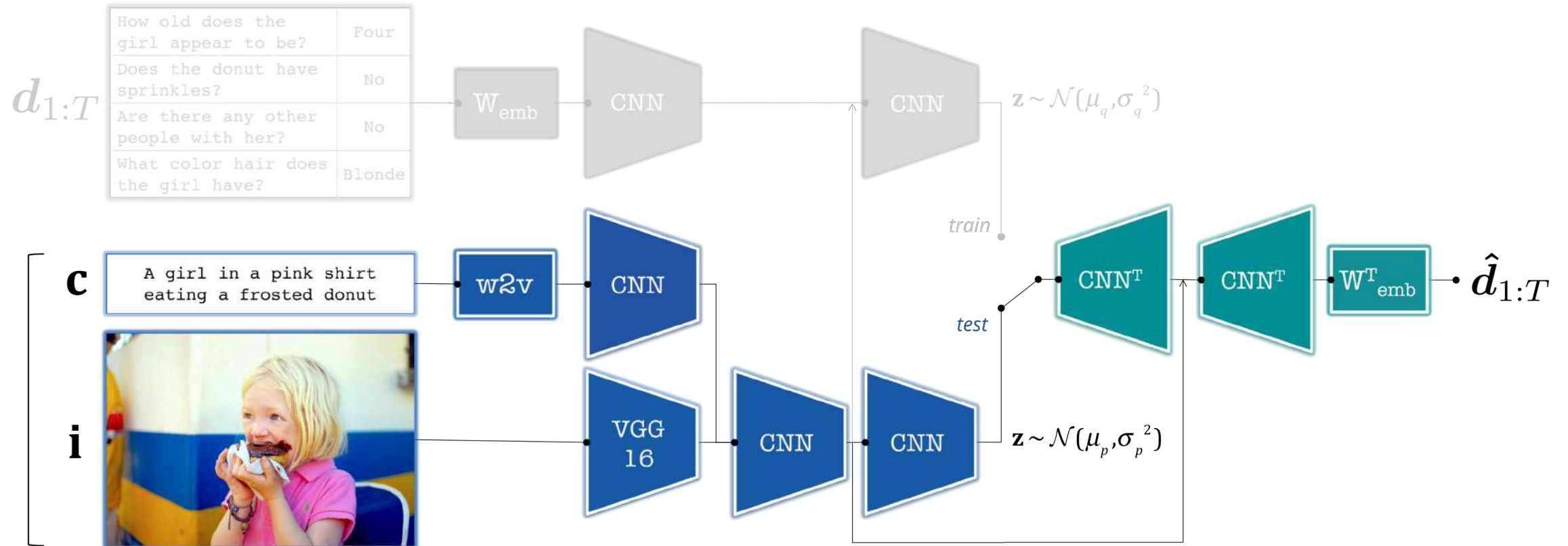




# AN AUTO-REGRESSIVE (AR) LAYER



# 2VD EVALUATION



## 2VD EVALUATION (I)

|                            | History? | NLL (↓)      |
|----------------------------|----------|--------------|
| <b>Ours</b>                | ∅        | 31.18        |
|                            | ✓        | 25.40        |
|                            |          |              |
| <b>Ours<sub>AR8</sub></b>  | ∅        | 28.81        |
|                            | ✓        | 26.60        |
|                            |          |              |
| <b>Ours<sub>AR10</sub></b> | ∅        | 28.49        |
|                            | ✓        | <b>24.93</b> |

∅ - generate whole dialogue block

✓ - iteratively condition on previously generated history



## ONE-WAY VISUAL DIALOGUE (IVD)

How many people  
are there?

it's crowded

Where is the  
pavement?

on the left

Where are the  
balloons?

in the air

A street with people and balloons

## TWO-WAY VISUAL DIALOGUE (2VD)

How many people  
are there?

it's crowded

Where is the  
pavement?

on the left

Where are the  
balloons?

in the air

A street with people and balloons

# LATENT SPACE DISPERSION ( $\text{SIM}_{\mathcal{U}}$ )

ground-truth dialogue  $d_{1:T}$

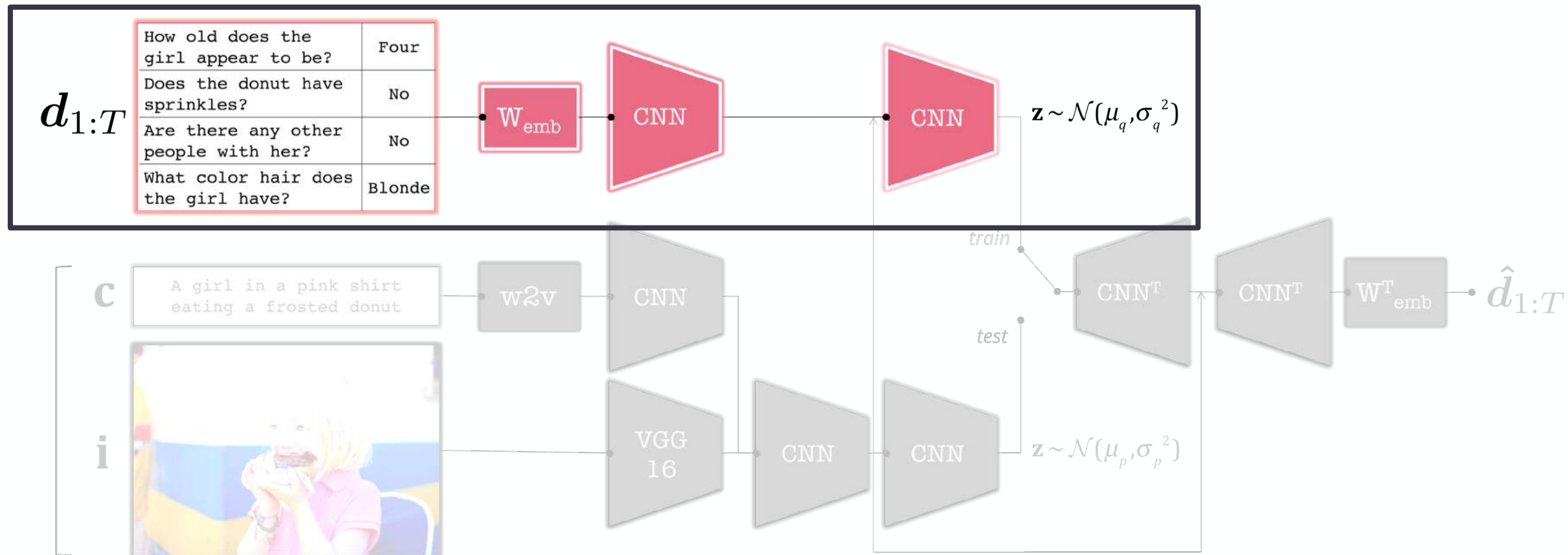
|                                      |        |
|--------------------------------------|--------|
| How old does the girl appear to be?  | Four   |
| Does the donut have sprinkles?       | No     |
| Are there any other people with her? | No     |
| What color hair does the girl have?  | Blonde |

generated dialogue  $\hat{d}_{1:T}$

|                           |              |
|---------------------------|--------------|
| What is the girl wearing? | A pink shirt |
| Is she eating?            | Yes          |
| Is she young or old?      | Young        |
| What type of donut is it? | Chocolate    |

$z$

# LATENT SPACE DISPERSION ( $\text{SIM}_{\mathcal{G}}$ )



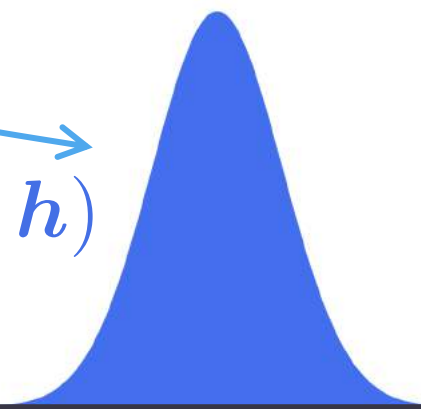


# LATENT SPACE DISPERSION ( $\text{SIM}_{\mathcal{U}}$ )

ground-truth dialogue  $d_{1:T}$

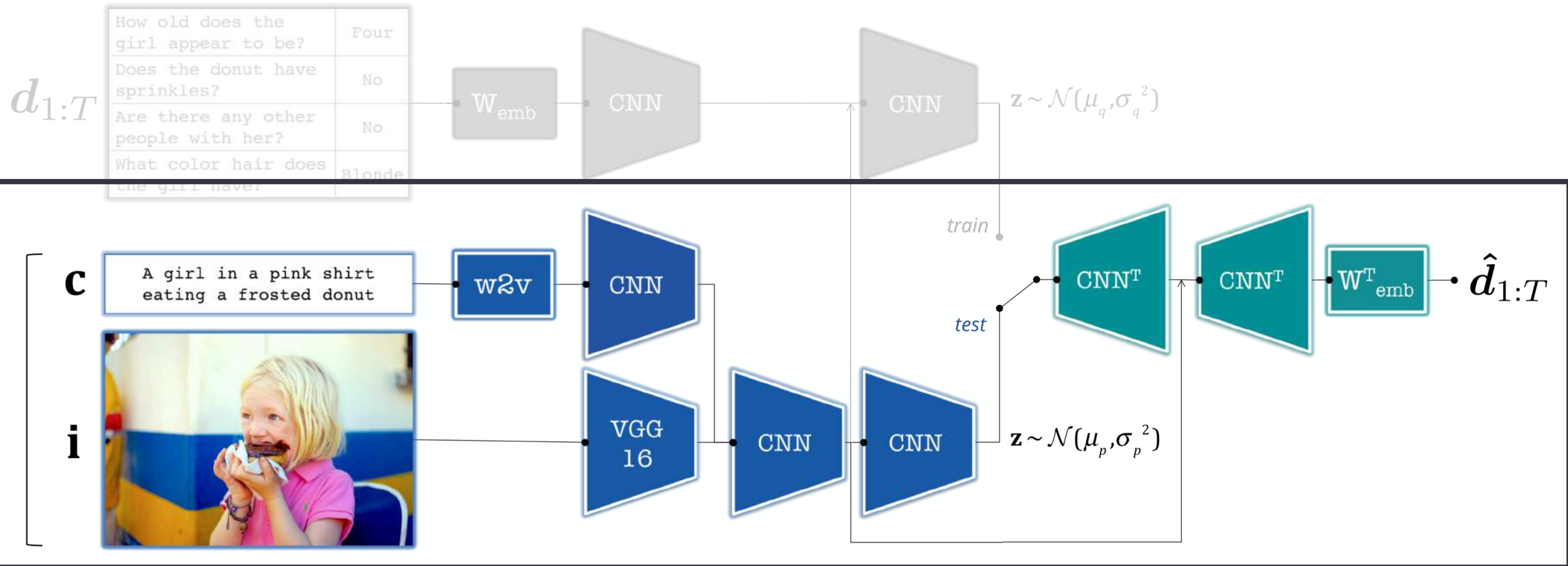
|                                      |        |
|--------------------------------------|--------|
| How old does the girl appear to be?  | Four   |
| Does the donut have sprinkles?       | No     |
| Are there any other people with her? | No     |
| What color hair does the girl have?  | Blonde |

$q_{\phi}(z \mid d_{1:T}, i, c, h)$

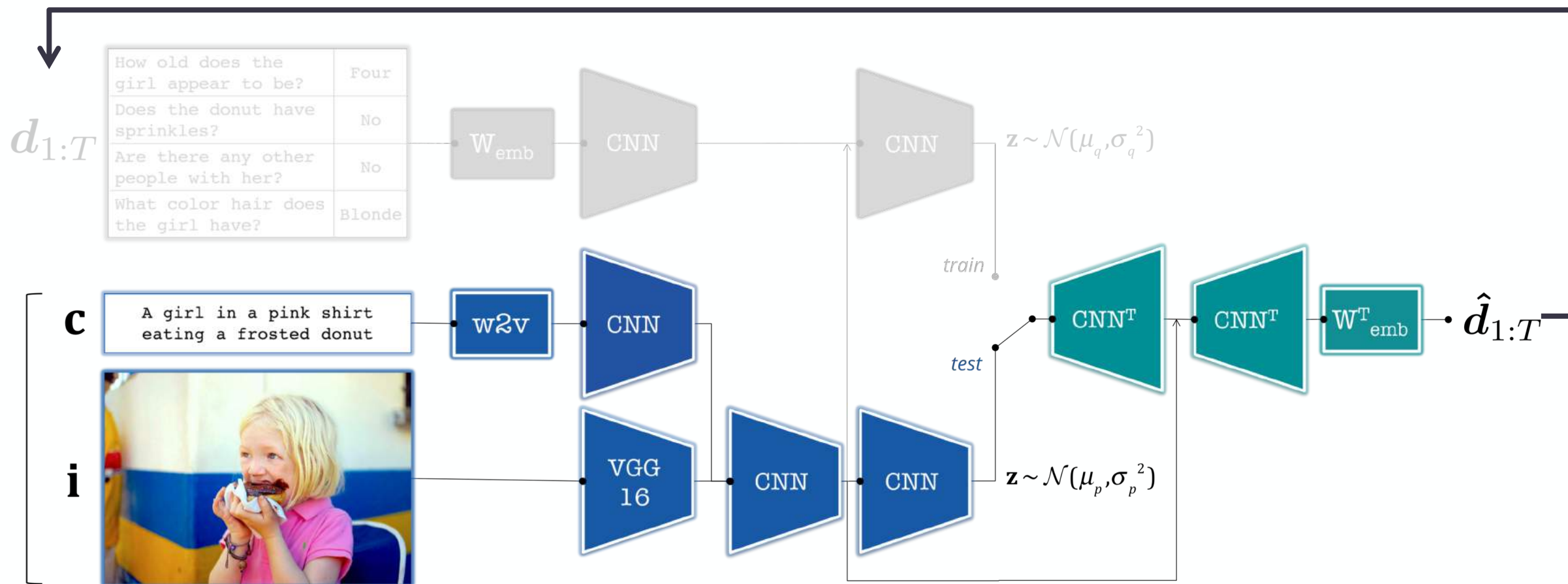


$z$

# LATENT SPACE DISPERSION ( $\text{SIM}_{\mathcal{G}}$ )

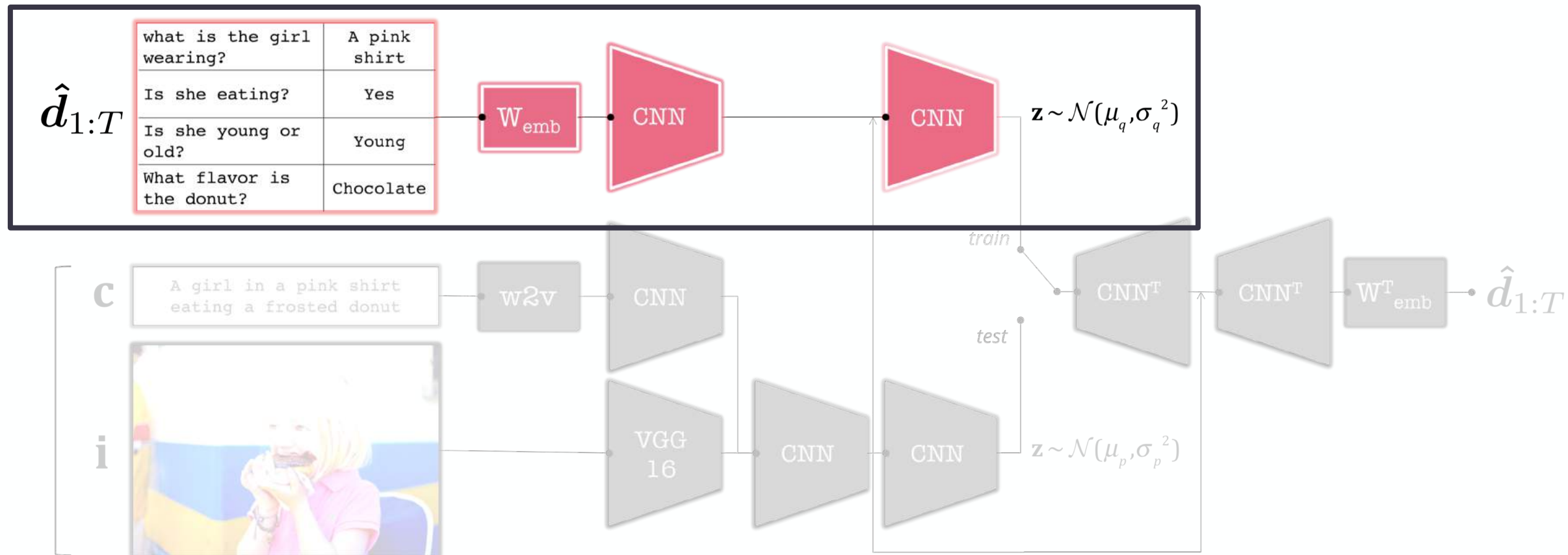


# LATENT SPACE DISPERSION ( $\text{SIM}_{\mathcal{U}}$ )





# LATENT SPACE DISPERSION (SIM<sub>Q</sub>)



# LATENT SPACE DISPERSION ( $\text{SIM}_{\mathcal{G}}$ )

ground-truth dialogue  $d_{1:T}$

|                                      |        |
|--------------------------------------|--------|
| How old does the girl appear to be?  | Four   |
| Does the donut have sprinkles?       | No     |
| Are there any other people with her? | No     |
| What color hair does the girl have?  | Blonde |

generated dialogue  $\hat{d}_{1:T}$

|                           |              |
|---------------------------|--------------|
| What is the girl wearing? | A pink shirt |
| Is she eating?            | Yes          |
| Is she young or old?      | Young        |
| What type of donut is it? | Chocolate    |

$$q_{\phi}(z \mid d_{1:T}, i, c, h)$$

$$q_{\phi}(z \mid \hat{d}_{1:T}, i, c, h)$$

$z$

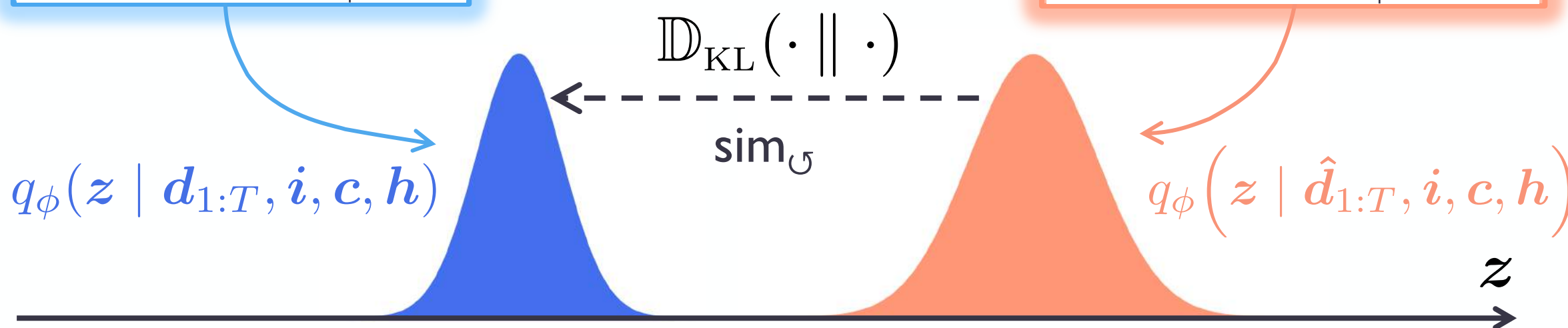
# LATENT SPACE DISPERSION ( $\text{SIM}_{\mathcal{U}}$ )

ground-truth dialogue  $d_{1:T}$

|                                      |        |
|--------------------------------------|--------|
| How old does the girl appear to be?  | Four   |
| Does the donut have sprinkles?       | No     |
| Are there any other people with her? | No     |
| What color hair does the girl have?  | Blonde |

generated dialogue  $\hat{d}_{1:T}$

|                           |              |
|---------------------------|--------------|
| What is the girl wearing? | A pink shirt |
| Is she eating?            | Yes          |
| Is she young or old?      | Young        |
| What type of donut is it? | Chocolate    |





# QUESTION RELEVANCY ( $\text{SIM}_{\text{cQ}}$ )

caption  $c$

A girl in a pink shirt  
eating a frosted donut

generated questions from  $\hat{d}_{1:T}$

{  
what is the girl wearing?  
is she eating?  
is she young or old?  
}

word2vec embedding  
space

# QUESTION RELEVANCY ( $\text{SIM}_{\text{cQ}}$ )

caption  $c$

A girl in a pink shirt  
eating a frosted donut

average  
word2vec

•  
 $c$

generated questions from  $\hat{d}_{1:T}$

{ what is the girl wearing?  
is she eating?  
is she young or old? }

word2vec embedding  
space

# QUESTION RELEVANCY ( $\text{SIM}_{\text{cQ}}$ )

caption  $c$

A girl in a pink shirt  
eating a frosted donut

average  
word2vec

•  
 $c$

generated questions from  $\hat{d}_{1:T}$

{  
what is the girl wearing?  
is she eating?  
is she young or old?  
}

•  
 $\hat{q}_1$

word2vec embedding  
space



# QUESTION RELEVANCY ( $\text{SIM}_{\text{cQ}}$ )

caption  $c$

A girl in a pink shirt  
eating a frosted donut

average  
word2vec

$\bullet$   
 $c$

generated questions from  $\hat{d}_{1:T}$

{ what is the girl wearing?  
**is she eating?**  
is she young or old? }

$\bullet$   $\hat{q}_2$

$\bullet$   $\hat{q}_1$

word2vec embedding  
space

# QUESTION RELEVANCY ( $\text{SIM}_{\text{cQ}}$ )

caption  $c$

A girl in a pink shirt  
eating a frosted donut

average  
word2vec

$\bullet c$

generated questions from  $\hat{d}_{1:T}$

{  
what is the girl wearing?  
is she eating?  
is she young or old?  
}

$\bullet \hat{q}_3$

$\bullet \hat{q}_2$

$\bullet \hat{q}_1$

word2vec embedding  
space

# QUESTION RELEVANCY ( $\text{SIM}_{\text{cQ}}$ )

caption  $c$

A girl in a pink shirt  
eating a frosted donut

average  
word2vec

$c$

cosine distance

$\hat{q}_3$

$\hat{q}_2$

$\hat{q}_1$

word2vec embedding  
space

generated questions from  $\hat{d}_{1:T}$

{  
what is the girl wearing?  
is she eating?  
is she young or old?  
}



## 2VD EVALUATION (II)

|                            | History? | NLL (↓)      | sim <sub>U</sub> (↓) | sim <sub>cq</sub> (↑) |
|----------------------------|----------|--------------|----------------------|-----------------------|
| <b>Ours</b>                | ∅        | 31.18        | 14.20                | <b>0.4931</b>         |
|                            | ✓        | 25.40        | <b>1.86</b>          | 0.4091                |
|                            |          |              |                      |                       |
| <b>Ours<sub>AR8</sub></b>  | ∅        | 28.81        | 31.50                | 0.4878                |
|                            | ✓        | 26.60        | 2.39                 | 0.3884                |
|                            |          |              |                      |                       |
| <b>Ours<sub>AR10</sub></b> | ∅        | 28.49        | 44.34                | 0.4927                |
|                            | ✓        | <b>24.93</b> | 2.35                 | 0.4101                |

∅ - generate whole dialogue block

✓ - iteratively condition on previously generated history

## 2VD EVALUATION (III)



Man and a boy playing  
ball in the grass

|                            |       |
|----------------------------|-------|
| Is the picture in color?   | Yes   |
| Is the photo in?           | Yes   |
| What color is the frisbee? | White |
| Can there?                 | No    |
| Is the grass visible?      | Yes   |
| Is the sunny?              | Yes   |
| Is it sunny?               | Yes   |
| Is there other people?     | No    |
| What color is the frisbee? | White |
| Is the man green?          | Yes   |

|                                 |       |
|---------------------------------|-------|
| What color is the man's?        | White |
| What is the man man?            | I     |
| How old is the man?             | I     |
| Can you see any ball?           | Yes   |
| Is there any other people?      | No    |
| Does the man have a?            | No    |
| Is the man wearing a hat?       | No    |
| Does he man have a?             | No    |
| Are there see any other people? | No    |
| Can you see any sky?            | No    |

## CONCLUSION & FUTURE

→ **IVD & 2VD!**

- ✓ Generative visual dialogue
- ✓ Diverse questions & answers
- ✓ Novel evaluation metrics

- + Increasing role of visual cues
- + Quantifying diversity
- + Asking informative questions



# FLIPDIAL IVD WEB-DEMO



Hello!  
I'm a  
visual  
chat-bot. Pick  
an image and  
quiz me!

enter your question...



[www.robots.ox.ac.uk/~daniela/research/flipdial/1vd\\_demo](http://www.robots.ox.ac.uk/~daniela/research/flipdial/1vd_demo)



# FLIPDIAL: FIND US AT POSTER #C18



✓ Generative visual dialogue

✓ Diverse questions & answers

✓ Novel evaluation metrics

+ Increasing role of visual cues

+ Quantifying diversity

+ Asking informative questions



Daniela Massiceti



N. Siddharth



Puneet Dokania



Philip H.S. Torr