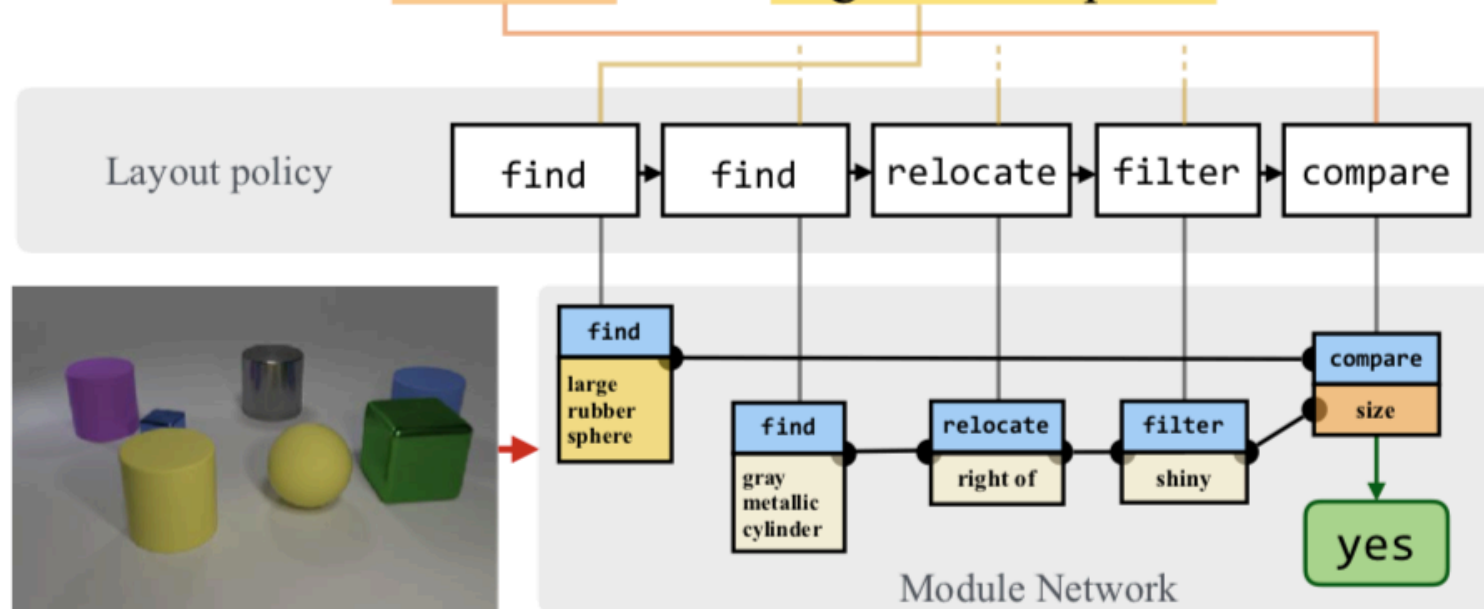


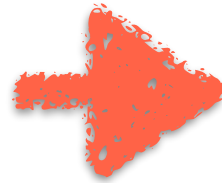
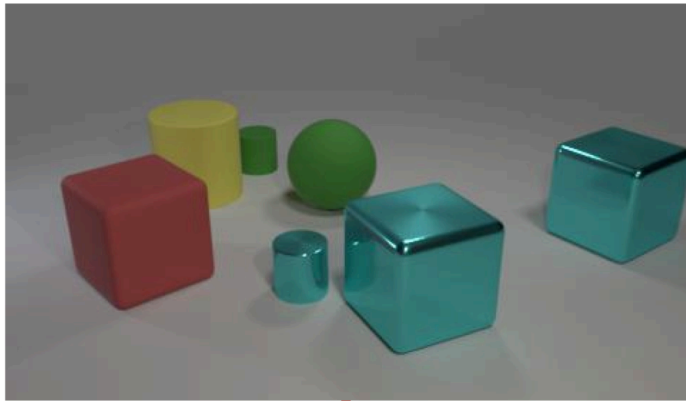
# Learning to Reason: End-to-End Module Networks for Visual Question Answering

There is a shiny object that is right of the gray metallic cylinder; does it have the same size as the large rubber sphere?

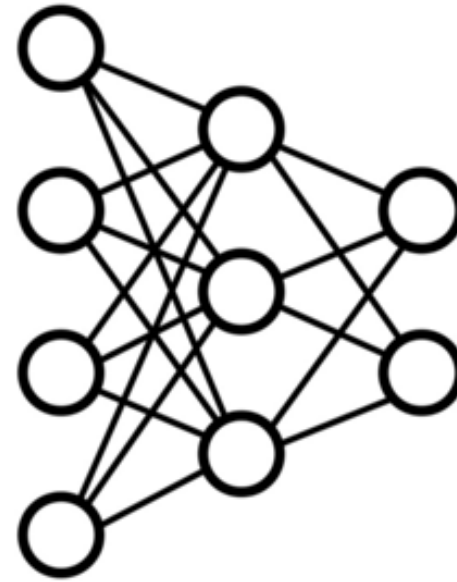


문 태 동

# Motivation



Monolithic Neural Network



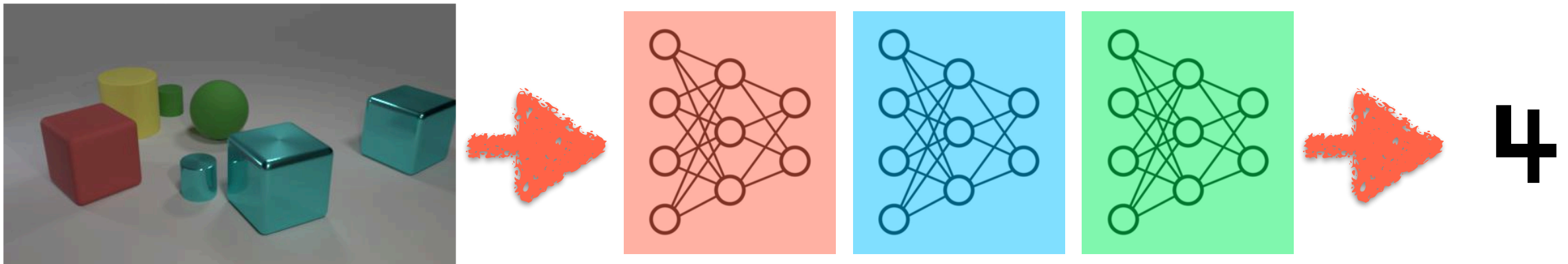
3

실제 답은 4

How many other things are of the same size as the green matte ball?

# Motivation

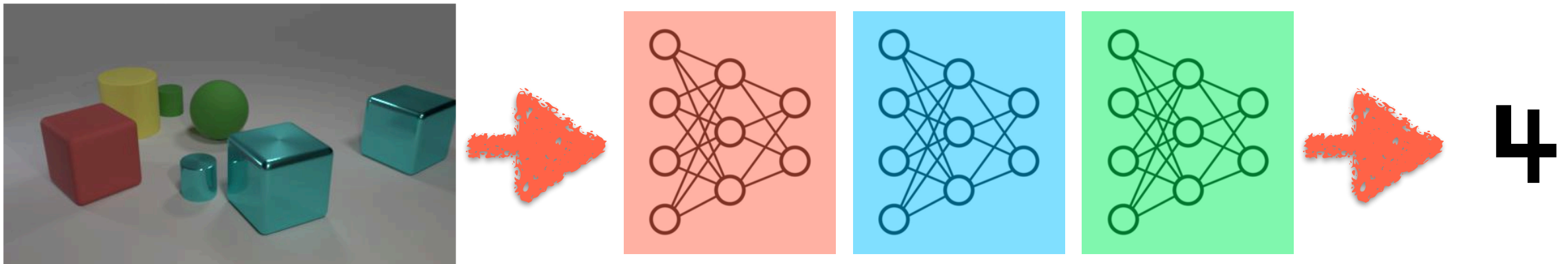
## Neural Module Networks



How many other things are of the same size as the green matte ball?

# Motivation

## Neural Module Networks



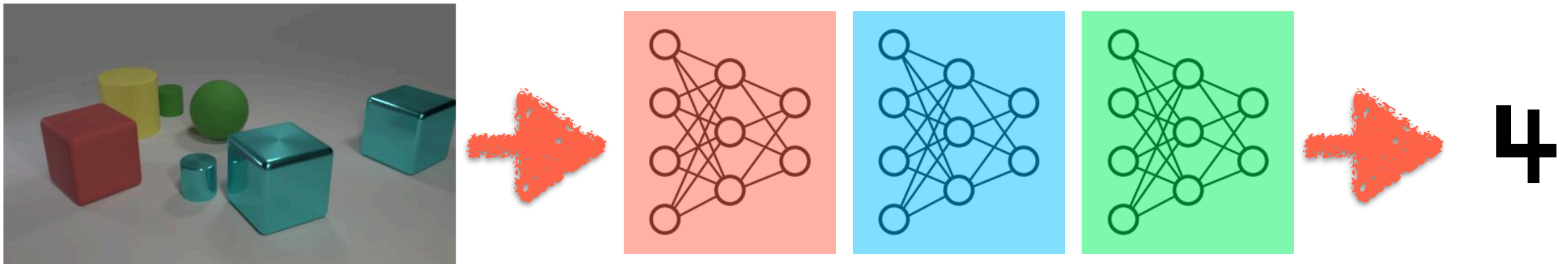
How many other things are of the same size as the green matte ball?

질문에 따라 Neural Module로 Network를 구성하자  
그러면 어떻게 구성할까?

policy:  $p(\text{Layout} \mid \text{Question})$ 로 Layout 예측

# Motivation

Neural Module Networks

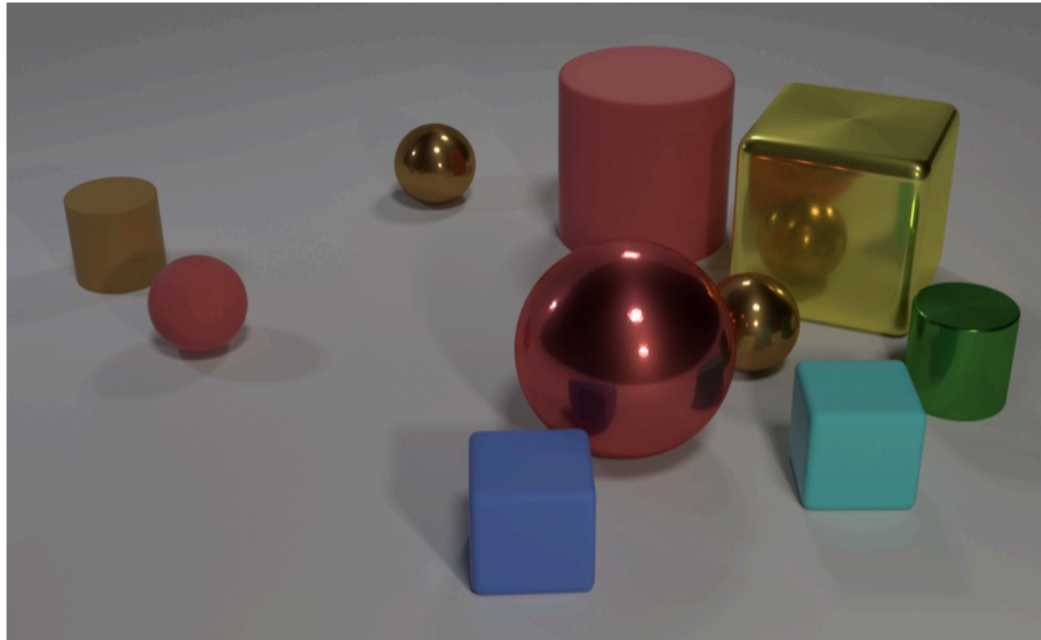


How many other things are of the same size as the green matte ball?

End-to-end neural network: Layout Policy + NMN 함께 학습

# Dataset: Clevr

Questions in CLEVR test various aspects of visual reasoning including **attribute identification**, **counting**, **comparison**, **spatial relationships**, and **logical operations**.



이미지: 100,000

질문: 1,000,000

Train(70%)/Val(15%)/Test(15%)

**Q:** Are there an **equal number** of **large things** and **metal spheres**?

**Q:** **What size** is the **cylinder that is left of** the **brown metal** thing **that is left of** the **big sphere**?

**Q:** There is a **sphere** with the **same size as** the **metal cube**; is it **made of the same material as** the **small red sphere**?

**Q:** **How many** objects are **either small cylinders** or **red** things?

# Dataset: Clevr

jupyter answers\_clevr.txt✓

File Edit View Language

```
1 0
2 1
3 10
4 2
5 3
6 4
7 5
8 6
9 7
10 8
11 9
12 blue
13 brown
14 cube
15 cyan
16 cylinder
17 gray
18 green
19 large
20 metal
21 no
22 purple
23 red
24 rubber
25 small
26 sphere
27 yellow
28 yes
29
```

jupyter vocabulary\_clevr.txt✓

File Edit View Language

```
49 matte
50 metal
51 metallic
52 more
53 number
54 object
55 objects
56 of
57 on
58 or
59 other
60 purple
61 red
62 right
63 rubber
64 same
65 shape
66 shiny
67 side
68 size
69 small
70 sphere
71 spheres
72 than
73 that
74 the
75 there
76 thing
77 things
78 tiny
79 to
80 visible
81 what
82 yellow
83
```

jupyter vocabulary\_layout.txt✓

File Edit View Language

```
1 _Scene
2 _Find
3 _Filter
4 _FindSameProperty
5 _Transform
6 _And
7 _Or
8 _Count
9 _Exist
10 _EqualNum
11 _MoreNum
12 _LessNum
13 _SameProperty
14 _Describe
15 <eos>
16 |
```

질문을 해봐요!

# Neural Modules

```

class FindModule(nn.Module):
    ...

    Mapping image_feat_grid X text_param -> att.grid
    (N,D_image,H,W) X (N,1,D_text) --> [N,1,H,W]
    ...

    def __init__(self, image_dim, text_dim, map_dim):
        super(FindModule, self).__init__()
        self.map_dim = map_dim
        self.conv1 = nn.Conv2d(image_dim, map_dim, kernel_size=1)
        self.conv2 = nn.Conv2d(map_dim, 1, kernel_size=1)
        self.textfc = nn.Linear(text_dim, map_dim)

    def forward(self, input_image_feat, input_text, input_image_attention1=None, input_image_attention2=None):
        image_mapped = self.conv1(input_image_feat) #(N, map_dim, H, W)
        text_mapped = self.textfc(input_text).view(-1, self.map_dim, 1, 1).expand_as(image_mapped)
        elmtwise_mult = image_mapped * text_mapped
        elmtwise_mult = F.normalize(elmtwise_mult, p=2, dim=1) #(N, map_dim, H, W)
        att_grid = self.conv2(elmtwise_mult) #(N, 1, H, W)
        return att_grid

```

Module name	Att-inputs	Features	Output	Implementation details
find	(none)	$x_{vis}, x_{txt}$	att	$a_{out} = \text{conv}_2(\text{conv}_1(x_{vis}) \odot W x_{txt})$
relocate	$a$	$x_{vis}, x_{txt}$	att	$a_{out} = \text{conv}_2(\text{conv}_1(x_{vis}) \odot W_1 \text{sum}(a \odot x_{vis}) \odot W_2 x_{txt})$
and	$a_1, a_2$	(none)	att	$a_{out} = \text{minimum}(a_1, a_2)$
or	$a_1, a_2$	(none)	att	$a_{out} = \text{maximum}(a_1, a_2)$
filter	$a$	$x_{vis}, x_{txt}$	att	$a_{out} = \text{and}(a, \text{find}[x_{vis}, x_{txt}]()), i.e. \text{reusing find and and}$
[exist, count]	$a$	(none)	ans	$y = W^T \text{vec}(a)$
describe	$a$	$x_{vis}, x_{txt}$	ans	$y = W_1^T (W_2 \text{sum}(a \odot x_{vis}) \odot W_3 x_{txt})$
[eq-count, more, less]	$a_1, a_2$	(none)	ans	$y = W_1^T \text{vec}(a_1) + W_2^T \text{vec}(a_2)$
compare	$a_1, a_2$	$x_{vis}, x_{txt}$	ans	$y = W_1^T (W_2 \text{sum}(a_1 \odot x_{vis}) \odot W_3 \text{sum}(a_2 \odot x_{vis}) \odot W_4 x_{txt})$



# Neural Modules

```
class AndModule(nn.Module):
    def __init__(self):
        super(AndModule, self).__init__()

    def forward(self, input_image_feat, input_text, input_image_attention1=None, input_image_attention2=None):
        return torch.max(input_image_attention1, input_image_attention2)

class OrModule(nn.Module):
    def __init__(self):
        super(OrModule, self).__init__()

    def forward(self, input_image_feat, input_text, input_image_attention1=None, input_image_attention2=None):
        return torch.min(input_image_attention1, input_image_attention2)
```

Module name	Att-inputs	Features	Output	Implementation details
find	(none)	$x_{vis}, x_{txt}$	att	$a_{out} = \text{conv}_2(\text{conv}_1(x_{vis}) \odot W x_{txt})$
relocate	$a$	$x_{vis}, x_{txt}$	att	$a_{out} = \text{conv}_2(\text{conv}_1(x_{vis}) \odot W_1 \text{sum}(a \odot x_{vis}) \odot W_2 x_{txt})$
and	$a_1, a_2$	(none)	att	$a_{out} = \text{minimum}(a_1, a_2)$
or	$a_1, a_2$	(none)	att	$a_{out} = \text{maximum}(a_1, a_2)$
filter	$a$	$x_{vis}, x_{txt}$	att	$a_{out} = \text{and}(a, \text{find}[x_{vis}, x_{txt}]()), i.e. \text{reusing find and and}$
[exist, count]	$a$	(none)	ans	$y = W^T \text{vec}(a)$
describe	$a$	$x_{vis}, x_{txt}$	ans	$y = W_1^T (W_2 \text{sum}(a \odot x_{vis}) \odot W_3 x_{txt})$
[eq-count, more, less]	$a_1, a_2$	(none)	ans	$y = W_1^T \text{vec}(a_1) + W_2^T \text{vec}(a_2)$
compare	$a_1, a_2$	$x_{vis}, x_{txt}$	ans	$y = W_1^T (W_2 \text{sum}(a_1 \odot x_{vis}) \odot W_3 \text{sum}(a_2 \odot x_{vis}) \odot W_4 x_{txt})$

# Neural Modules

```
class FilterModule(nn.Module):
    def __init__(self, findModule, andModule):
        super(FilterModule, self).__init__()
        self.andModule = andModule
        self.findModule = findModule

    def forward(self, input_image_feat, input_text, input_image_attention1=None, input_image_attention2=None):
        find_result = self.findModule(input_image_feat, input_text, input_image_attention1, input_image_attention2)
        att_grid = self.andModule(input_image_feat, input_text, input_image_attention1, find_result)
        return att_grid
```

Module name	Att-inputs	Features	Output	Implementation details
find	(none)	$x_{vis}, x_{txt}$	att	$a_{out} = \text{conv}_2(\text{conv}_1(x_{vis}) \odot W x_{txt})$
relocate	$a$	$x_{vis}, x_{txt}$	att	$a_{out} = \text{conv}_2(\text{conv}_1(x_{vis}) \odot W_1 \text{sum}(a \odot x_{vis}) \odot W_2 x_{txt})$
and	$a_1, a_2$	(none)	att	$a_{out} = \text{minimum}(a_1, a_2)$
or	$a_1, a_2$	(none)	att	$a_{out} = \text{maximum}(a_1, a_2)$
filter	$a$	$x_{vis}, x_{txt}$	att	$a_{out} = \text{and}(a, \text{find}[x_{vis}, x_{txt}]()), i.e. \text{reusing find and and}$
[exist, count]	$a$	(none)	ans	$y = W^T \text{vec}(a)$
describe	$a$	$x_{vis}, x_{txt}$	ans	$y = W_1^T (W_2 \text{sum}(a \odot x_{vis}) \odot W_3 x_{txt})$
[eq-count, more, less]	$a_1, a_2$	(none)	ans	$y = W_1^T \text{vec}(a_1) + W_2^T \text{vec}(a_2)$
compare	$a_1, a_2$	$x_{vis}, x_{txt}$	ans	$y = W_1^T (W_2 \text{sum}(a_1 \odot x_{vis}) \odot W_3 \text{sum}(a_2 \odot x_{vis}) \odot W_4 x_{txt})$

# End-to-End Neural Module Network (N2NMN)

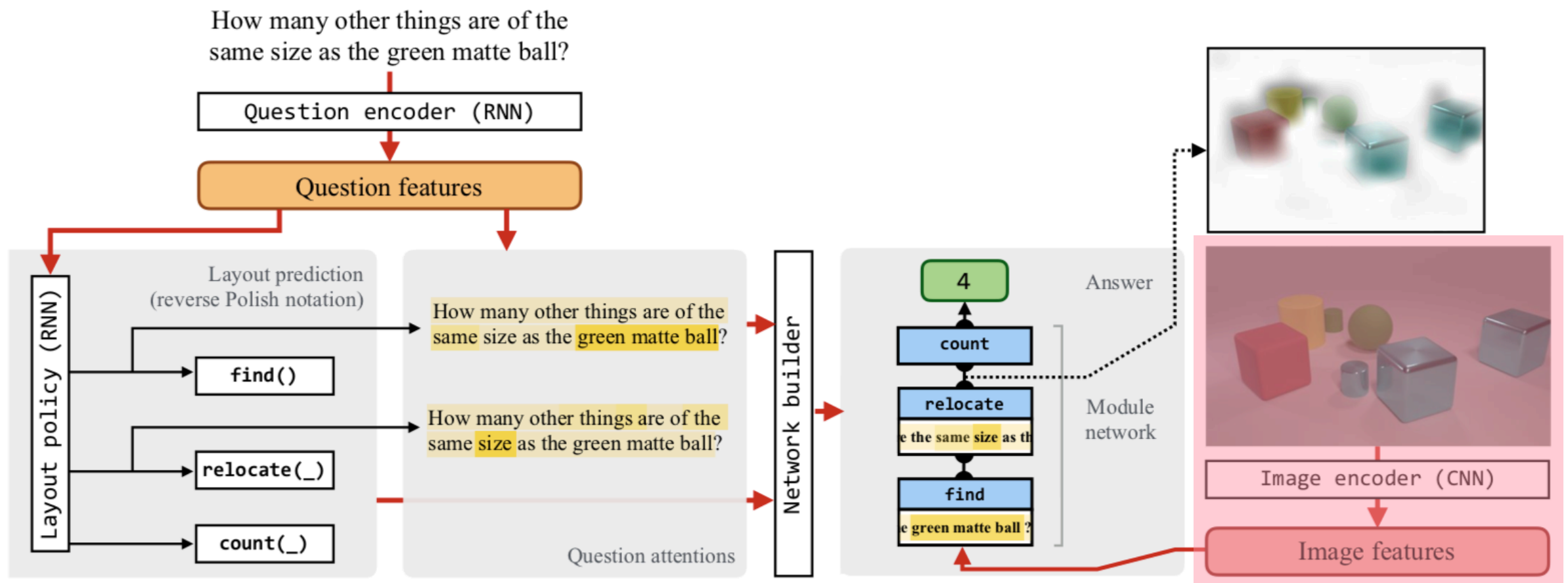


image feature extraction

# Image Feature Extraction

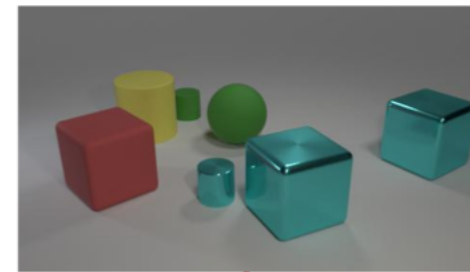


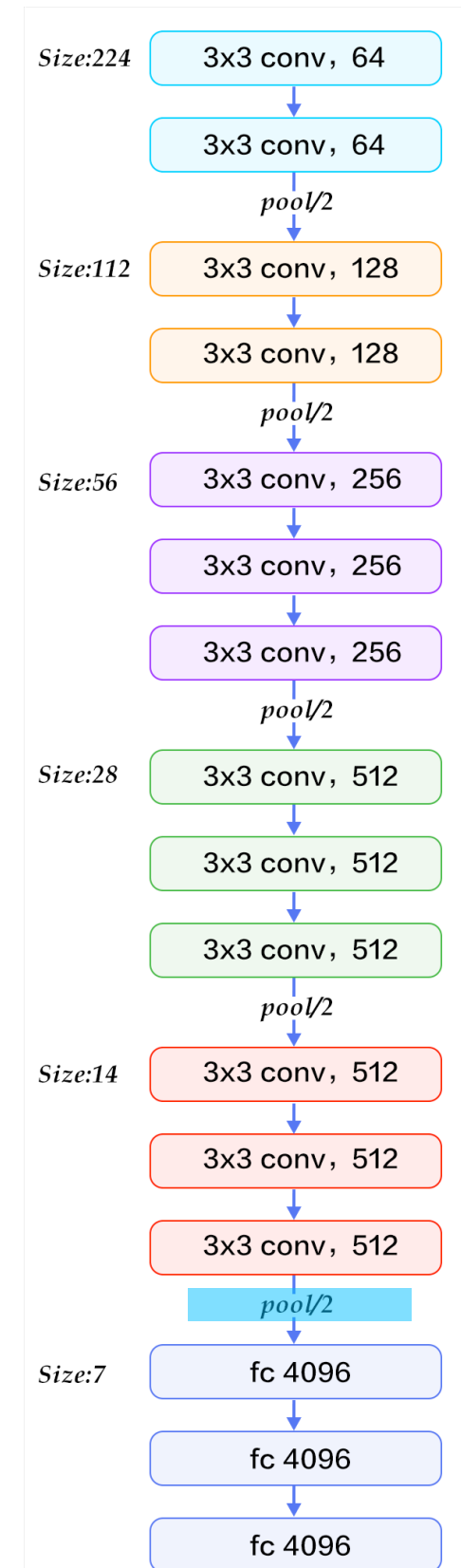
Image encoder (CNN)

Image features

## Image feature VGG-16

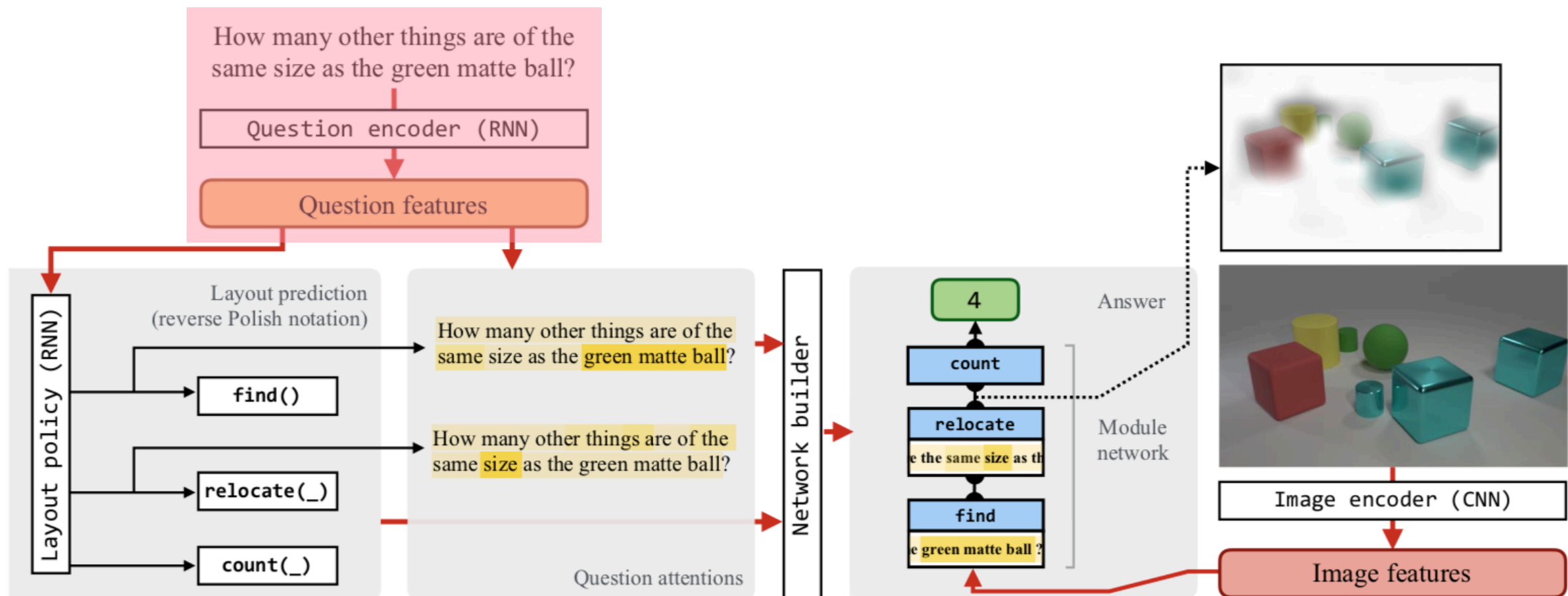
### TF code

```
24 # layer 2
25 conv2_1 = conv_relu('conv2_1', pool1,
26                     kernel_size=3, stride=1, output_dim=128)
27 conv2_2 = conv_relu('conv2_2', conv2_1,
28                     kernel_size=3, stride=1, output_dim=128)
29 pool2 = pool('pool2', conv2_2, kernel_size=2, stride=2)
30 # layer 3
31 conv3_1 = conv_relu('conv3_1', pool2,
32                     kernel_size=3, stride=1, output_dim=256)
33 conv3_2 = conv_relu('conv3_2', conv3_1,
34                     kernel_size=3, stride=1, output_dim=256)
35 conv3_3 = conv_relu('conv3_3', conv3_2,
36                     kernel_size=3, stride=1, output_dim=256)
37 pool3 = pool('pool3', conv3_3, kernel_size=2, stride=2)
38 # layer 4
39 conv4_1 = conv_relu('conv4_1', pool3,
40                     kernel_size=3, stride=1, output_dim=512)
41 conv4_2 = conv_relu('conv4_2', conv4_1,
42                     kernel_size=3, stride=1, output_dim=512)
43 conv4_3 = conv_relu('conv4_3', conv4_2,
44                     kernel_size=3, stride=1, output_dim=512)
45 pool4 = pool('pool4', conv4_3, kernel_size=2, stride=2)
46 # layer 5
47 conv5_1 = conv_relu('conv5_1', pool4,
48                     kernel_size=3, stride=1, output_dim=512)
49 conv5_2 = conv_relu('conv5_2', conv5_1,
50                     kernel_size=3, stride=1, output_dim=512)
51 conv5_3 = conv_relu('conv5_3', conv5_2,
52                     kernel_size=3, stride=1, output_dim=512)
53 pool5 = pool('pool5', conv5_3, kernel_size=2, stride=2)
54 return pool5
```

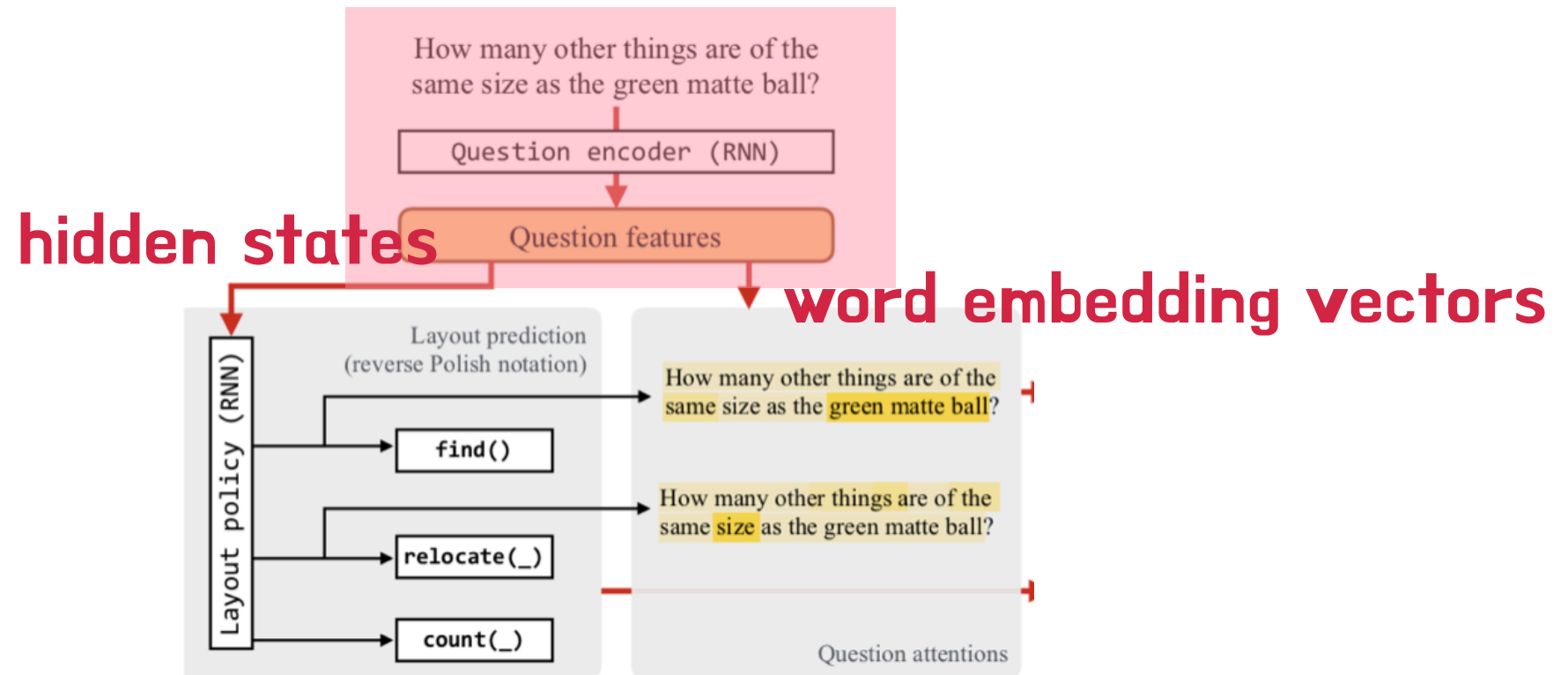


# End-to-End Neural Module Network (N2NMN)

## question feature extraction (encoding)



# Question Feature Extraction



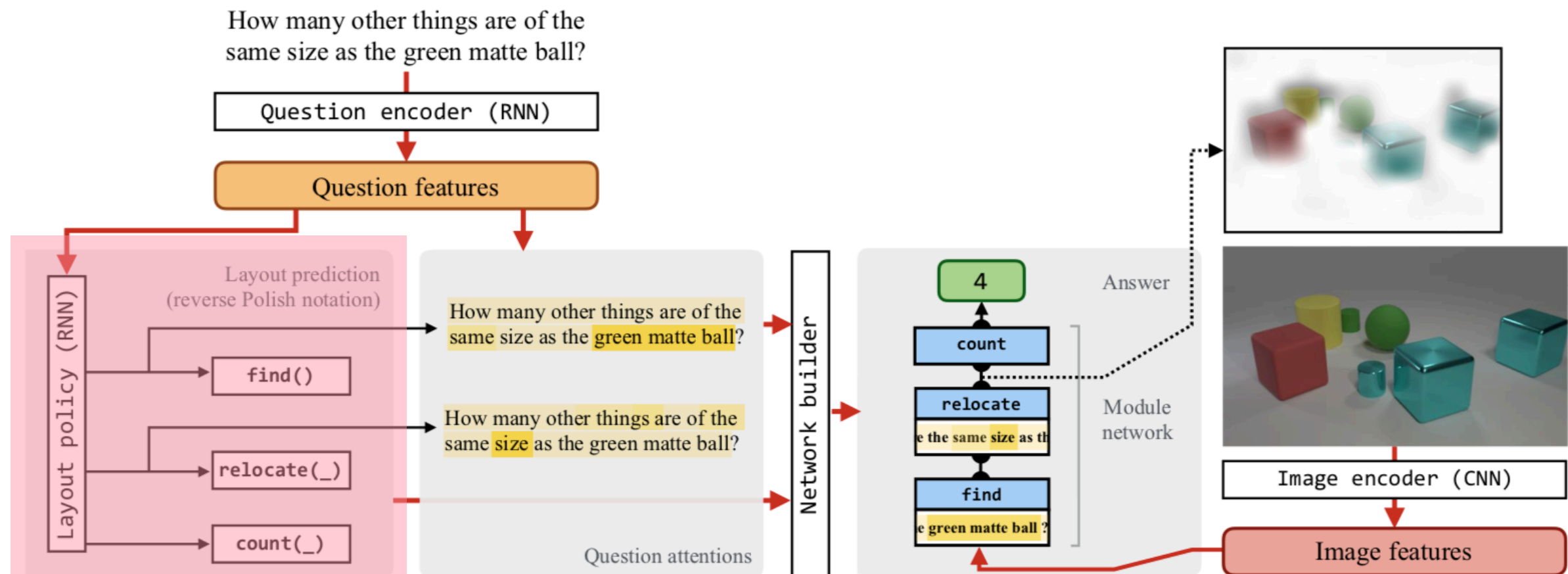
```
9 class EncoderRNN(nn.Module):
10     def __init__(self, input_size, hidden_size, input_encoding_size, num_layers=2):
11         super(EncoderRNN, self).__init__()
12         self.hidden_size = hidden_size
13         self.num_layers = num_layers
14
15         self.embedding = nn.Embedding(input_size, input_encoding_size)
16         self.lstm = nn.LSTM(input_encoding_size, hidden_size)
17
18     def forward(self, input_seqs, input_seq_lens, hidden):
19         embedded = self.embedding(input_seqs)
20         outputs, hidden = self.lstm(embedded)
21         return outputs, hidden, embedded
22
23     def initHidden(self, batch_size):
24         result = torch.zeros(self.num_layers, batch_size, self.hidden_size)
25         return result
```

## cf) LSTM 코드 구조

- outputs: 모든 time step에서 hidden states
- hidden: 마지막 time step에서 hidden state



# End-to-End Neural Module Network (N2NMN)



layout prediction (decoding)

# Layout expression

How many other things are of the same size as the green matte ball?

functional  
expression

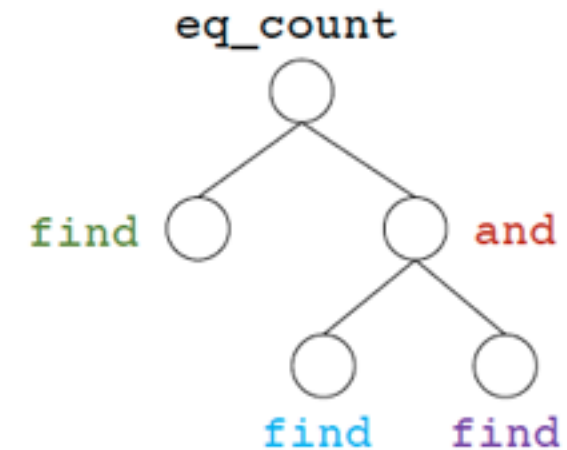


**count(relocate(find))**

layout  
expression

eq\_count(find(), and(find(), find()))

syntax tree



Reverse Polish  
Notation

[find, find, find, and, eq\_count]

NM Layout 표기법

NM을 network로 assemble하는 순서



따라서, Seq-to-seq network로 <eos>까지 NM token을 예측하여 assemble하자!

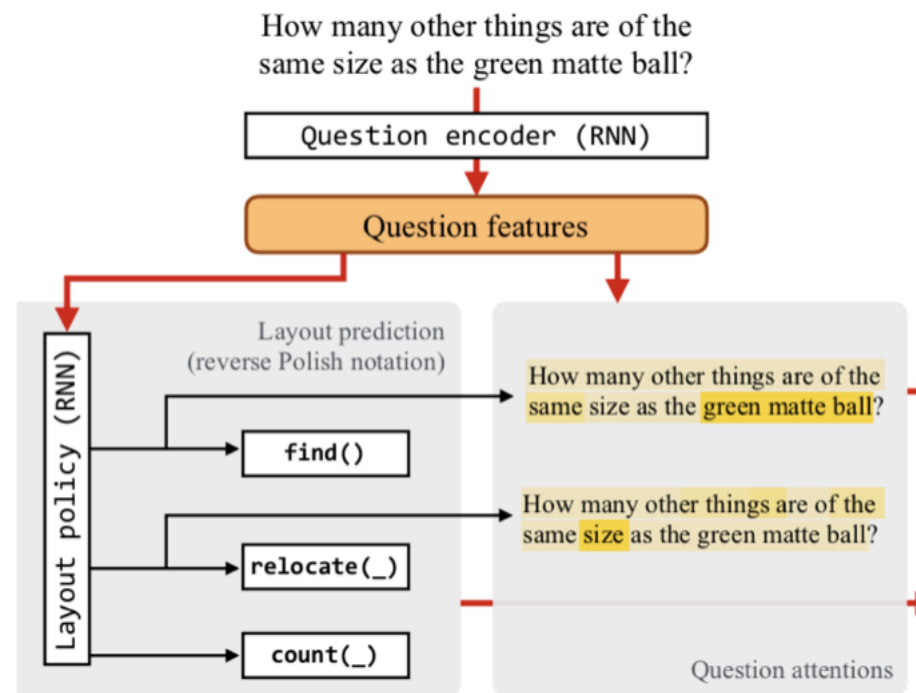
Translation  
English token -> French token



NM token 예측  
English token -> NM token

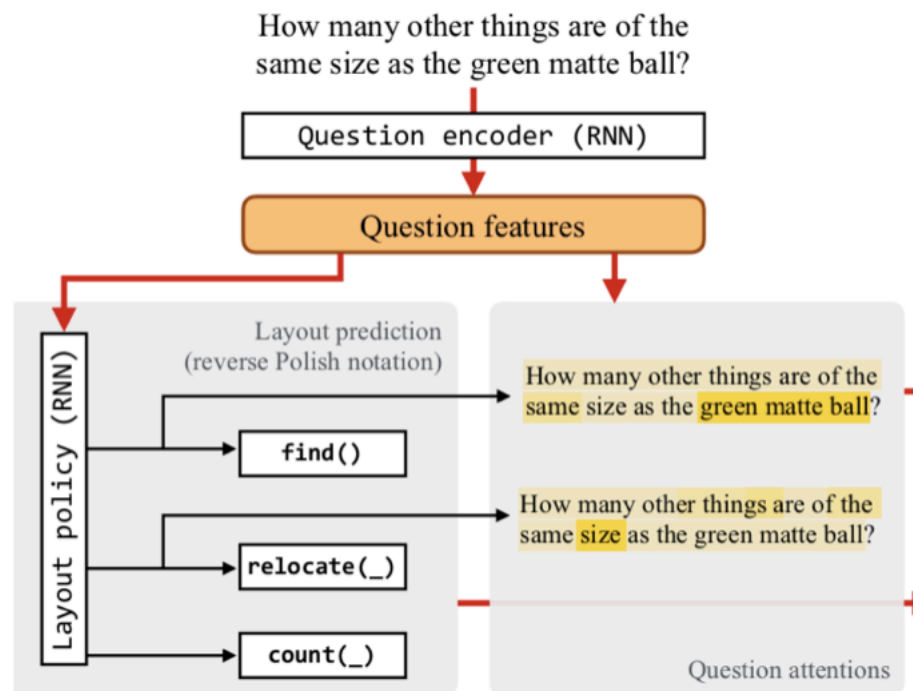


# End-to-End Neural Module Network (N2NMN)



Seq-to-seq net with attention 구조는  
코드와 함께 보며 칠판에^^

# End-to-End Neural Module Network (N2NMN)



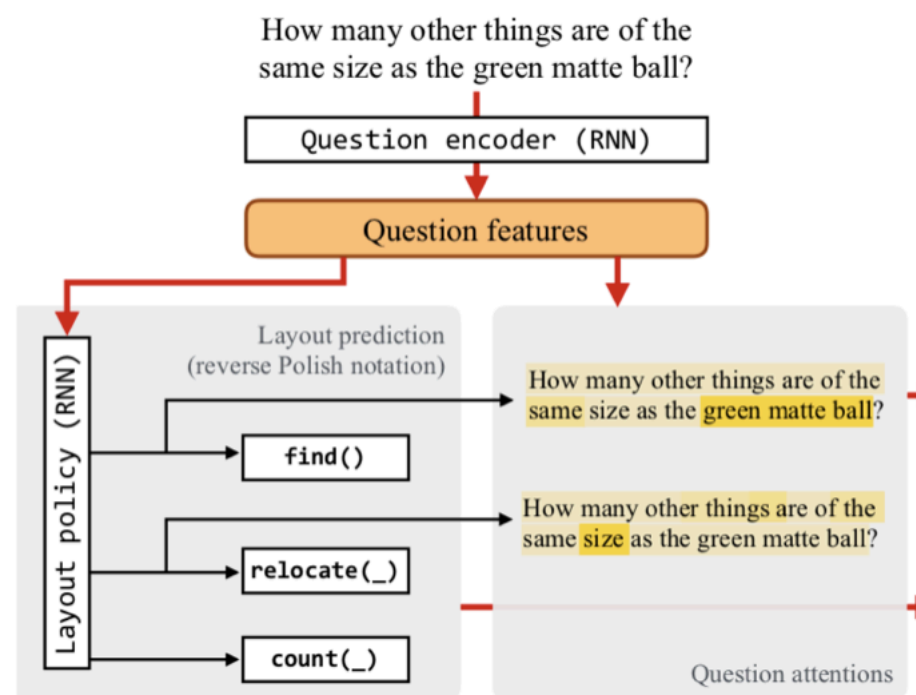
```
class LuongAttnDecoderRNN(nn.Module):
    def __init__(self, attn_model, embedding, hidden_size, output_size, n_layers=1,
                 dropout=0.1):
        super(LuongAttnDecoderRNN, self).__init__()

        # Keep for reference
        self.attn_model = attn_model
        self.hidden_size = hidden_size
        self.output_size = output_size
        self.n_layers = n_layers
        self.dropout = dropout

        # Define layers
        self.embedding = embedding
        self.embedding_dropout = nn.Dropout(dropout)
        self.gru = nn.GRU(hidden_size, hidden_size, n_layers, dropout=(0 if n_layers == 1
else dropout))
        self.concat = nn.Linear(hidden_size * 2, hidden_size)
        self.out = nn.Linear(hidden_size, output_size)

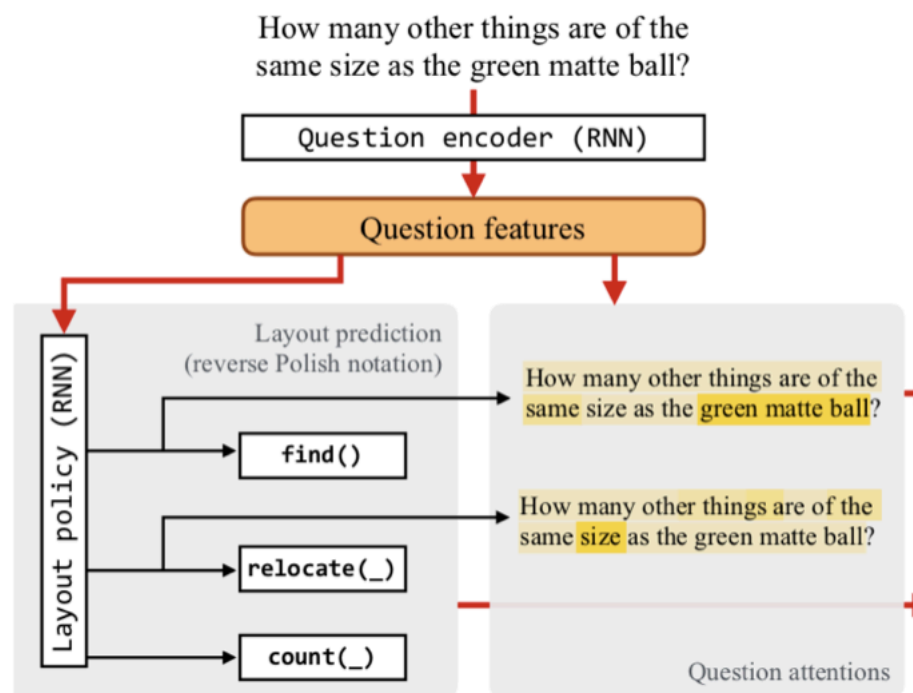
        self.attn = Attn(attn_model, hidden_size)
```

# End-to-End Neural Module Network (N2NMN)



```
def forward(self, input_step, last_hidden, encoder_outputs):  
    # Note: we run this one step (word) at a time  
    # Get embedding of current input word  
    embedded = self.embedding(input_step)  
    embedded = self.embedding_dropout(embedded)  
    # Forward through unidirectional GRU  
    rnn_output, hidden = self.gru(embedded, last_hidden)  
    # Calculate attention weights from the current GRU output  
    attn_weights = self.attn(rnn_output, encoder_outputs)  
    # Multiply attention weights to encoder outputs to get new "weighted sum" context  
    vector  
    context = attn_weights.bmm(encoder_outputs.transpose(0, 1))  
    # Concatenate weighted context vector and GRU output using Luong eq. 5  
    rnn_output = rnn_output.squeeze(0)  
    context = context.squeeze(1)  
    concat_input = torch.cat((rnn_output, context), 1)  
    concat_output = torch.tanh(self.concat(concat_input))  
    # Predict next word using Luong eq. 6  
    output = self.out(concat_output)  
    output = F.softmax(output, dim=1)  
    # Return output and final hidden state  
    return output, hidden
```

# End-to-End Neural Module Network (N2NMN)



```
# Luong attention layer
class Attn(torch.nn.Module):
    def __init__(self, method, hidden_size):
        super(Attn, self).__init__()
        self.method = method
        if self.method not in ['dot', 'general', 'concat']:
            raise ValueError(self.method, "is not an appropriate attention method.")
        self.hidden_size = hidden_size
        if self.method == 'general':
            self.attn = torch.nn.Linear(self.hidden_size, hidden_size)
        elif self.method == 'concat':
            self.attn = torch.nn.Linear(self.hidden_size * 2, hidden_size)
            self.v = torch.nn.Parameter(torch.FloatTensor(hidden_size))

    def dot_score(self, hidden, encoder_output):
        return torch.sum(hidden * encoder_output, dim=2)

    def general_score(self, hidden, encoder_output):
        energy = self.attn(encoder_output)
        return torch.sum(hidden * energy, dim=2)

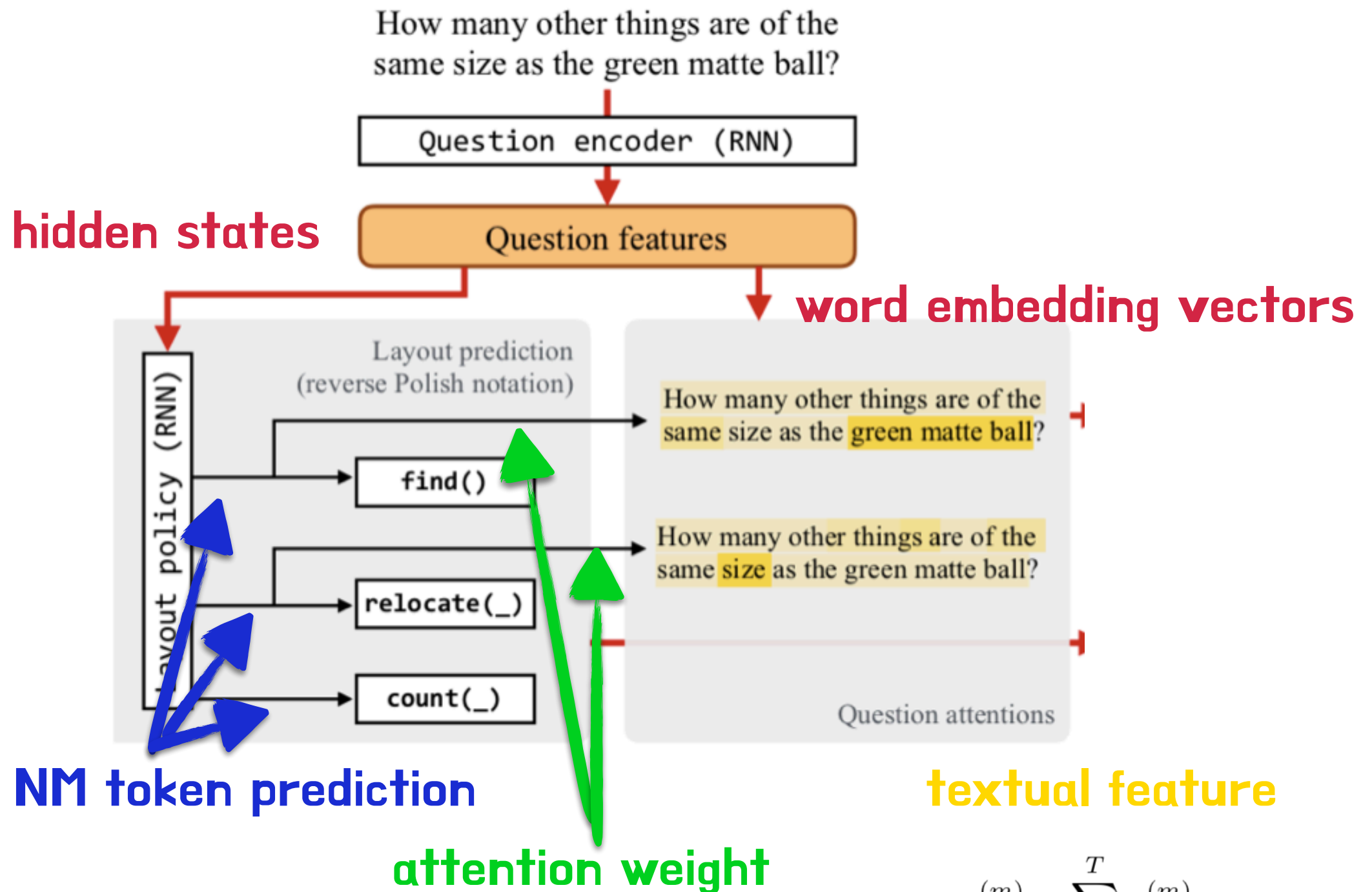
    def concat_score(self, hidden, encoder_output):
        energy = self.attn(torch.cat((hidden.expand(encoder_output.size(0), -1, -1),
encoder_output), 2)).tanh()
        return torch.sum(self.v * energy, dim=2)

    def forward(self, hidden, encoder_outputs):
        # Calculate the attention weights (energies) based on the given method
        if self.method == 'general':
            attn_energies = self.general_score(hidden, encoder_outputs)
        elif self.method == 'concat':
            attn_energies = self.concat_score(hidden, encoder_outputs)
        elif self.method == 'dot':
            attn_energies = self.dot_score(hidden, encoder_outputs)

        # Transpose max_length and batch_size dimensions
        attn_energies = attn_energies.t()

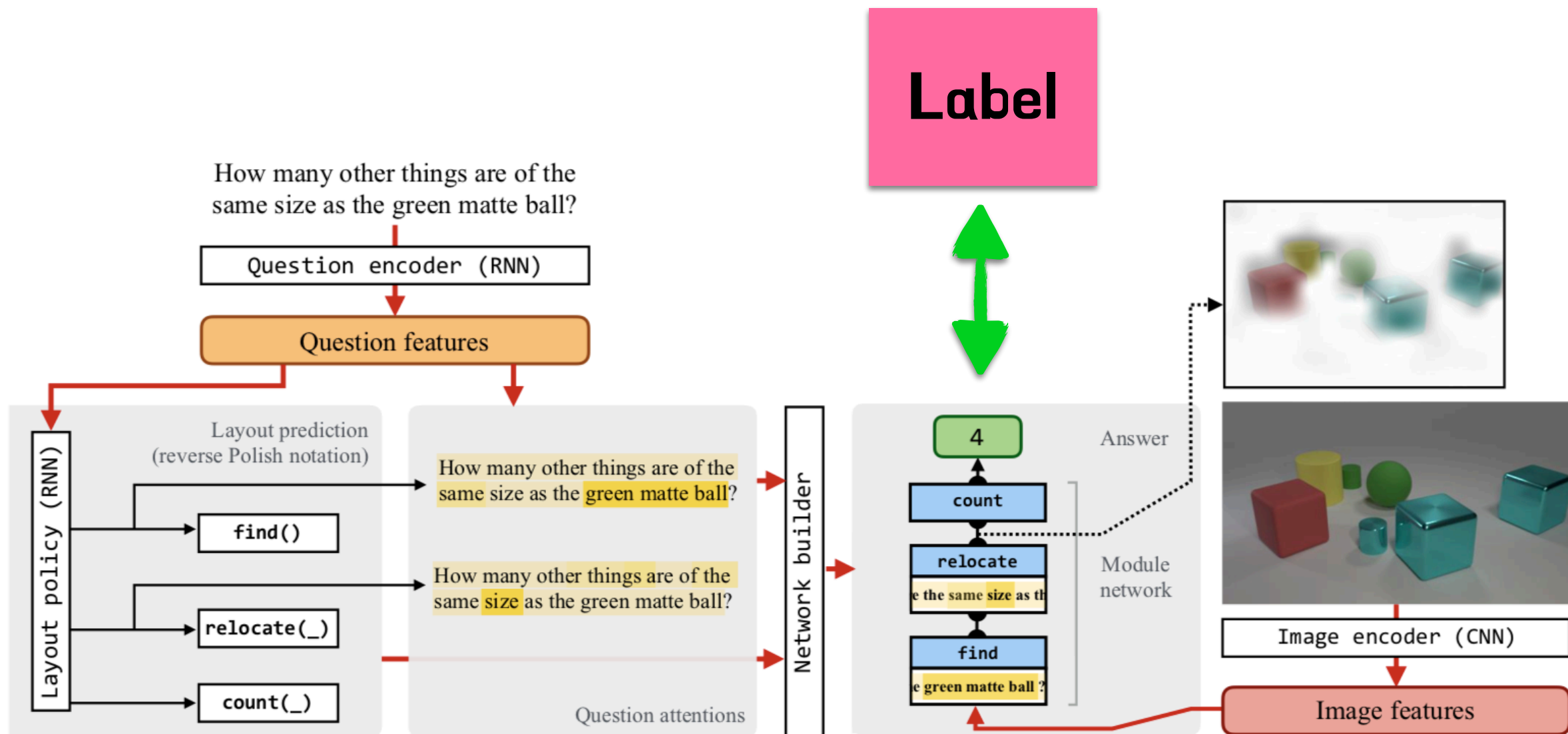
        # Return the softmax normalized probability scores (with added dimension)
        return F.softmax(attn_energies, dim=1).unsqueeze(1)
```

# NM token prediction + textual feature extraction



$$x_{txt}^{(m)} = \sum_{i=1}^T \alpha_i^{(m)} w_i$$

# End-to-end training



$$L(\theta) = E_{l \sim p(l|q; \theta)} [\tilde{L}(\theta, l; q, I)]$$

$$\nabla_{\theta} L \approx \frac{1}{M} \sum_{m=1}^M \left( \tilde{L}(\theta, l_m) \nabla_{\theta} \log p(l_m | q; \theta) + \nabla_{\theta} \tilde{L}(\theta, l_m) \right)$$

**끝!!**

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^\top \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^\top \tanh(\mathbf{W}_a [\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases}$$

