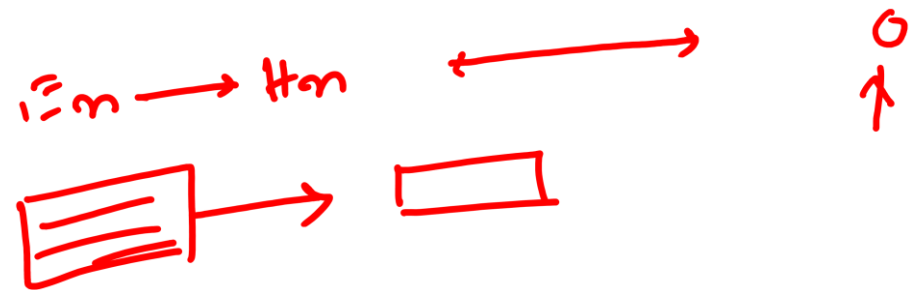# Sequence-to-Sequence Modeling

# Qwen3-Coder-Flash Unleashed

High-Performance Code Generation with Agent Integration!

Announced on August 1, 2025

Qwen3-Coder

> QWEN3-CODER-FLASH 🚀

Qwen3-Coder-30B-A3B-Instruct is a fine-tuned MoE (Mixture-of-Experts) variant in the Qwen3 model family: it has 30.5B total parameters, comprised of 128 experts, with only 8 experts (≈ 3.3B parameters) activated per inference - **making it highly efficient while retaining strong coding and reasoning performance.**
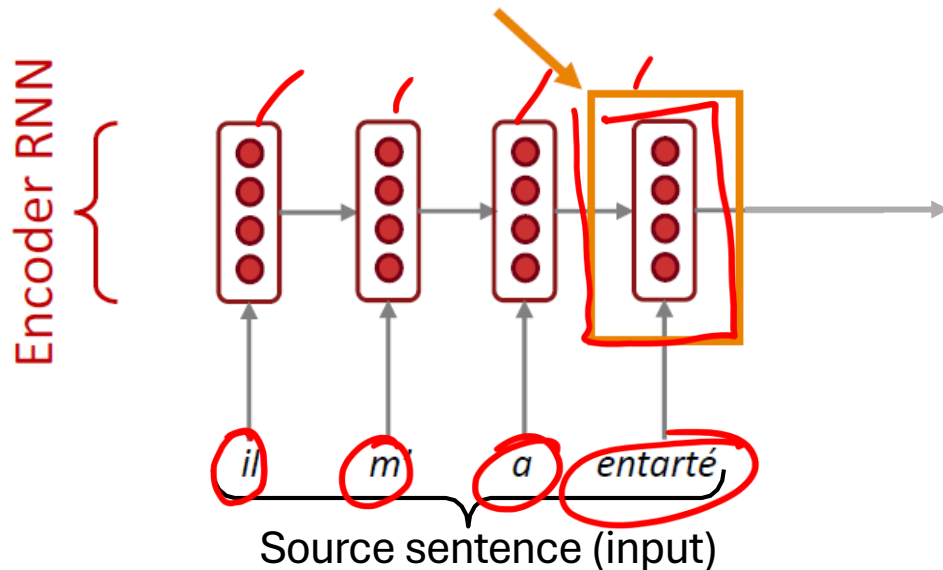
It shows significant performance among open models on Agentic Coding, Agentic Browser-Use, and other foundational coding tasks. It also features **Long-context Capabilities** with native support for **256K** tokens, extendable up to **1M** tokens using Yarn, optimized for repository-scale understanding.

| Benchmarks | Open Models | | | | Proprietary Models | |
|---|---|---|---|---|---|---|
| | Qwen3-Coder 30B-A3B-Instruct | Qwen3-Coder 480B-A35B-Instruct | Kimi-K2 Instruct | DeepSeek-V3 0324 | Claude Sonnet-4 | OpenAI GPT-4.1 |
| **Agentic Coding** | | | | | | |
| Terminal-Bench | 31.3 | 37.5 | 30.0 | 2.5 | 35.5 | 25.3 |
| SWE-bench Verified | | | | | | |
| w/ OpenHands, 500 turns | 51.6 | 69.6 | - | - | 70.4 | - |
| w/ OpenHands, 100 turns | 51.6 | 67.0 | 65.4 | 38.8 | 68.0 | 48.6 |
| w/ Private Scaffolding | - | - | 65.8 | - | 72.7 | 63.8 |
| SWE-bench Live | 20.7 | 26.3 | 22.3 | 13.0 | 27.7 | - |
| SWE-bench Multilingual | 34.7 | 54.7 | 47.3 | 13.0 | 53.3 | 31.5 |
| Multi-SWE-bench mini | 19.5 | 25.8 | 19.8 | 7.5 | 24.8 | - |
| Multi-SWE-bench flash | 19.3 | 27.0 | 20.7 | - | 25.0 | - |
| Aider-Polyglot | 33.3 | 61.8 | 60.0 | 56.9 | 56.4 | 52.4 |
| Spider2 | 21.4 | 31.1 | 25.2 | 17.7 | 31.1 | 25.6 |
| **Agentic Browser Use** | | | | | | |
| WebArena | 43.5 | 49.9 | 47.4 | 40.0 | 51.1 | 44.3 |
| Mind2Web | 51.0 | 55.8 | 42.7 | 36.0 | 47.4 | 49.6 |
| **Agentic Tool Use** | | | | | | |
| BFCL-v3 | 62.2 | 68.7 | 65.2 | 64.7 | 73.3 | 62.9 |
| TAU-Bench Retail | 68.7 | 77.5 | 70.7 | 59.1 | 80.5 | - |
| TAU-Bench Airline | 48.0 | 60.0 | 53.5 | 40.0 | 60.0 | - |

# Neural Machine Translation (NMT)

**The Sequence-to-Sequence Model**



Encoding of the source sentence.
Provides initial hidden state
for Decoder RNN.

Encoder RNN

il    m'    a    entarté

Source sentence (input)

**Encoder RNN produces an encoding of the source sentence.**

# Neural Machine Translation (NMT)

*Teacher-forcing*

**The Sequence-to-Sequence Model**

Encoding of the source sentence.
Provides initial hidden state
for Decoder RNN.

Target sentence (output)

he    hit    me    with    a    pie    <END>

<START>    he    hit    me    with    a    pie

Source sentence (input)

il    m'    a    entarté

Encoder RNN

Decoder RNN

argmax

**Decoder RNN** is a
**Language Model that
generates target
sentence,** *conditioned
on* *encoding*.

**Encoder RNN produces an encoding of the source sentence.**

**Note:** This diagram shows **test time**
behavior: decoder output is fed in as
next step's input

# Training an NMT System

Seq2seq is optimized as a **single system.** Backpropagation operates "*end-to-end*".

$$J = \frac{1}{T}\sum_{t=1}^{T} J_t$$

$= J_1 + J_2 + J_3 + J_4 + J_5 + J_6 + J_7 \Rightarrow J$

= negative log prob of "he"

= negative log prob of "with"

= negative log prob of <END>

$\hat{y}_1 \quad \hat{y}_2 \quad \hat{y}_3 \quad \hat{y}_4 \quad \hat{y}_5 \quad \hat{y}_6 \quad \hat{y}_7$

Encoder RNN

Decoder RNN

il    m'    a    entarté    <START>    he    hit    me    with    a    pie

Source sentence (from corpus)

Target sentence (from corpus)

# Issues With RNN

- Linear interaction distance
- Bottleneck problem
- Lack of parallelizability

**ATTENTION**
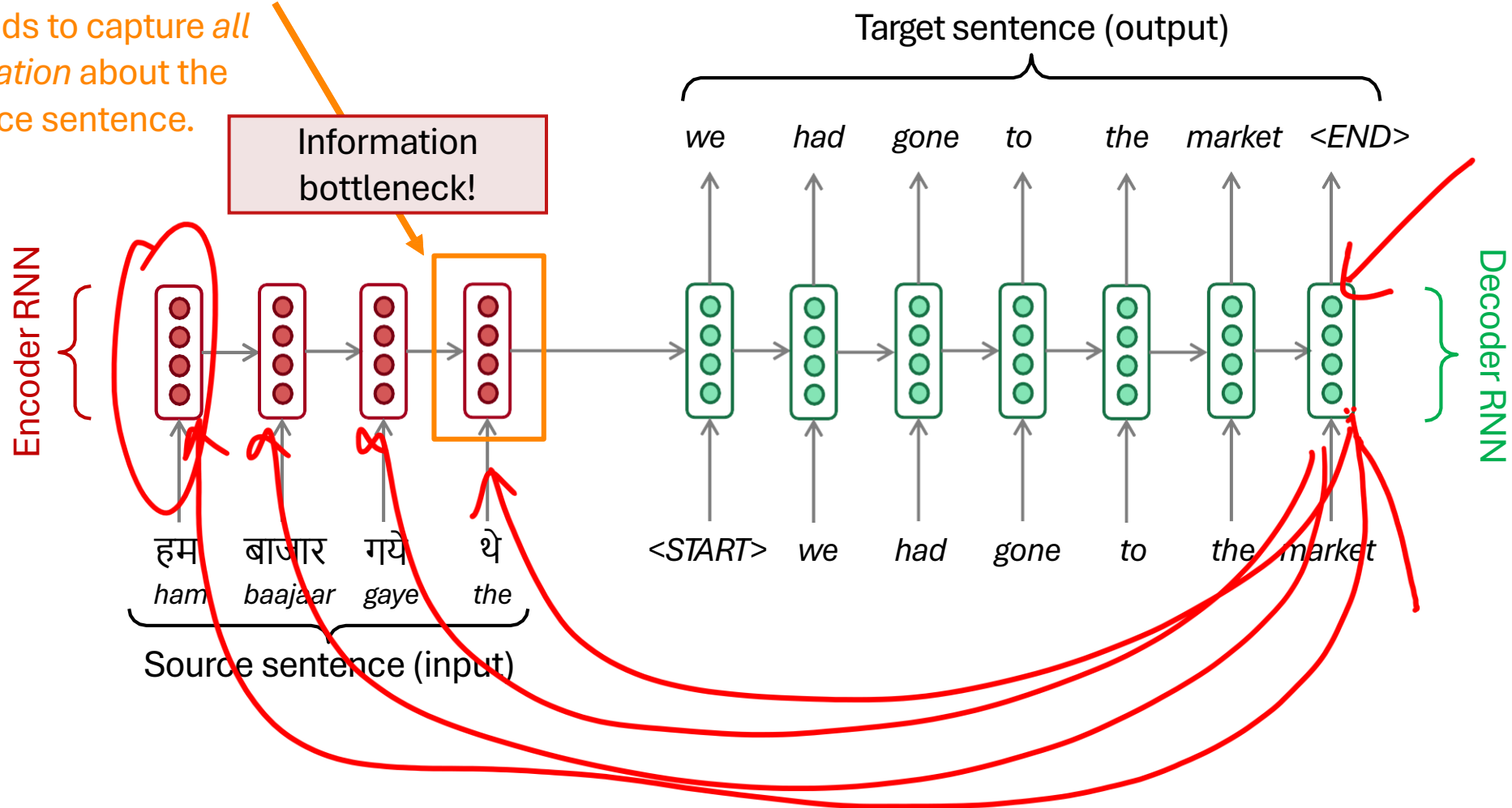
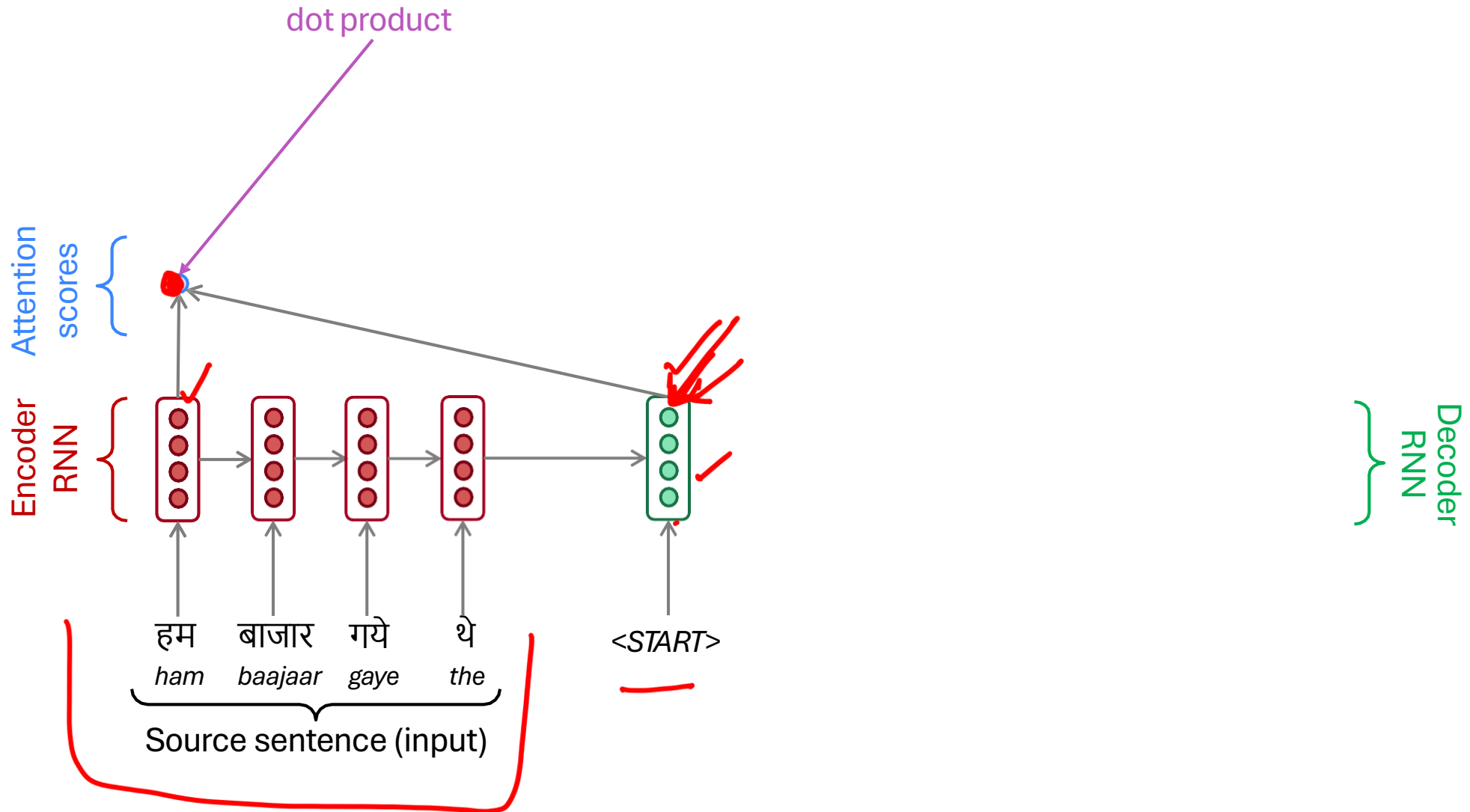# Sequence-to-Sequence: The Bottleneck Problem

Encoding of the source sentence

Target sentence (output)

Encoder RNN

Decoder RNN

*we    had    gone    to    the    market    <END>*

हम    बाजार    गये    थे
*ham    baajaar    gaye    the*

Source sentence (input)

*<START>    we    had    gone    to    the    market*

**Any problems with this architecture?**

# Sequence-to-Sequence: The Bottleneck Problem

Encoding of the source sentence
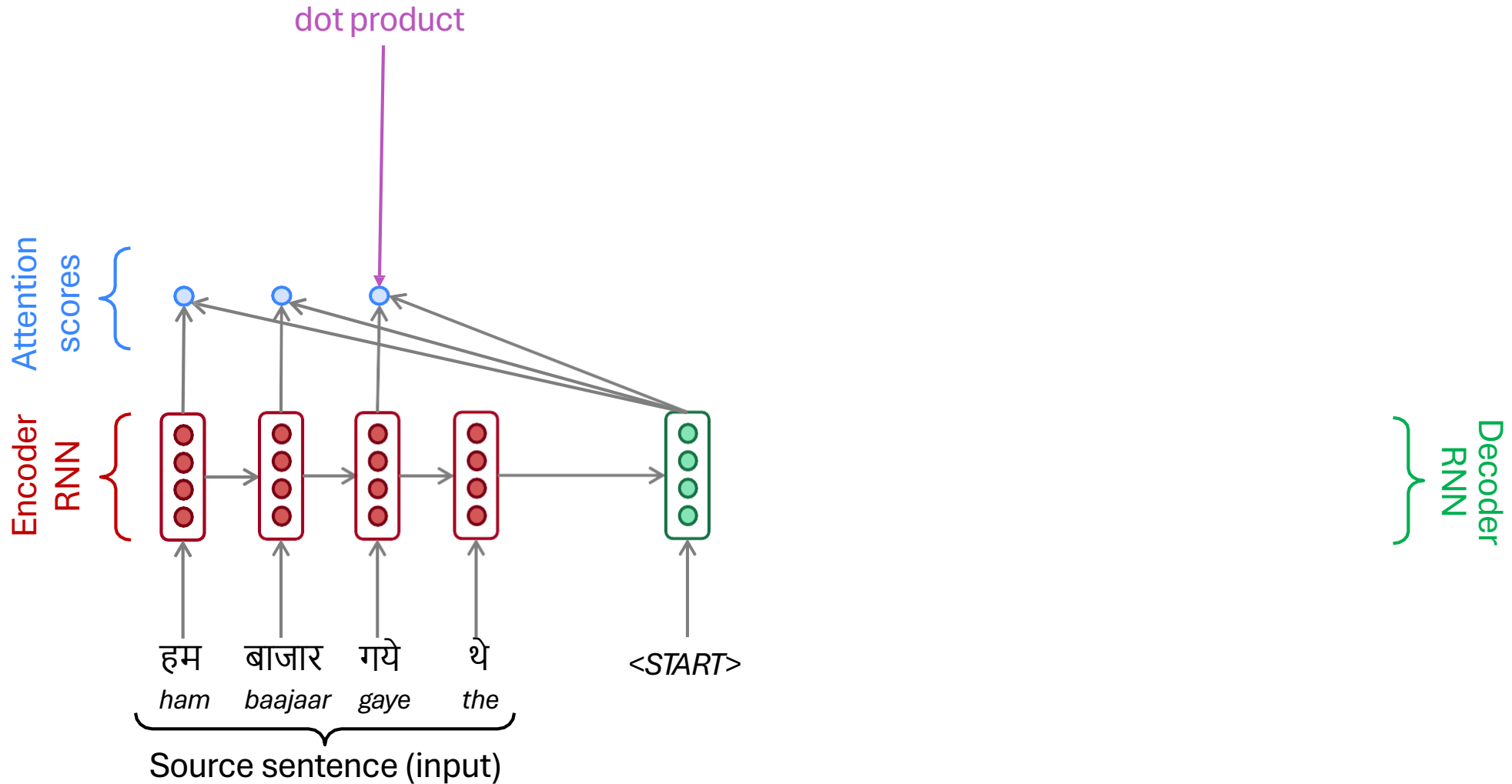This needs to capture *all information* about the source sentence.

Information bottleneck!

Target sentence (output)

Encoder RNN

Decoder RNN

we    had    gone    to    the    market    <END>

हम    बाजार    गये    थे
ham    baajaar    gaye    the

Source sentence (input)

<START>    we    had    gone    to    the    market

# Sequence-to-Sequence With Attention

dot product

Attention scores

Encoder RNN

Decoder RNN



हम
*ham*

बाजार
*baajaar*

गये
*gaye*

थे
*the*

<START>

Source sentence (input)

# Sequence-to-Sequence With Attention



dot product

Attention scores

Encoder RNN

Decoder RNN

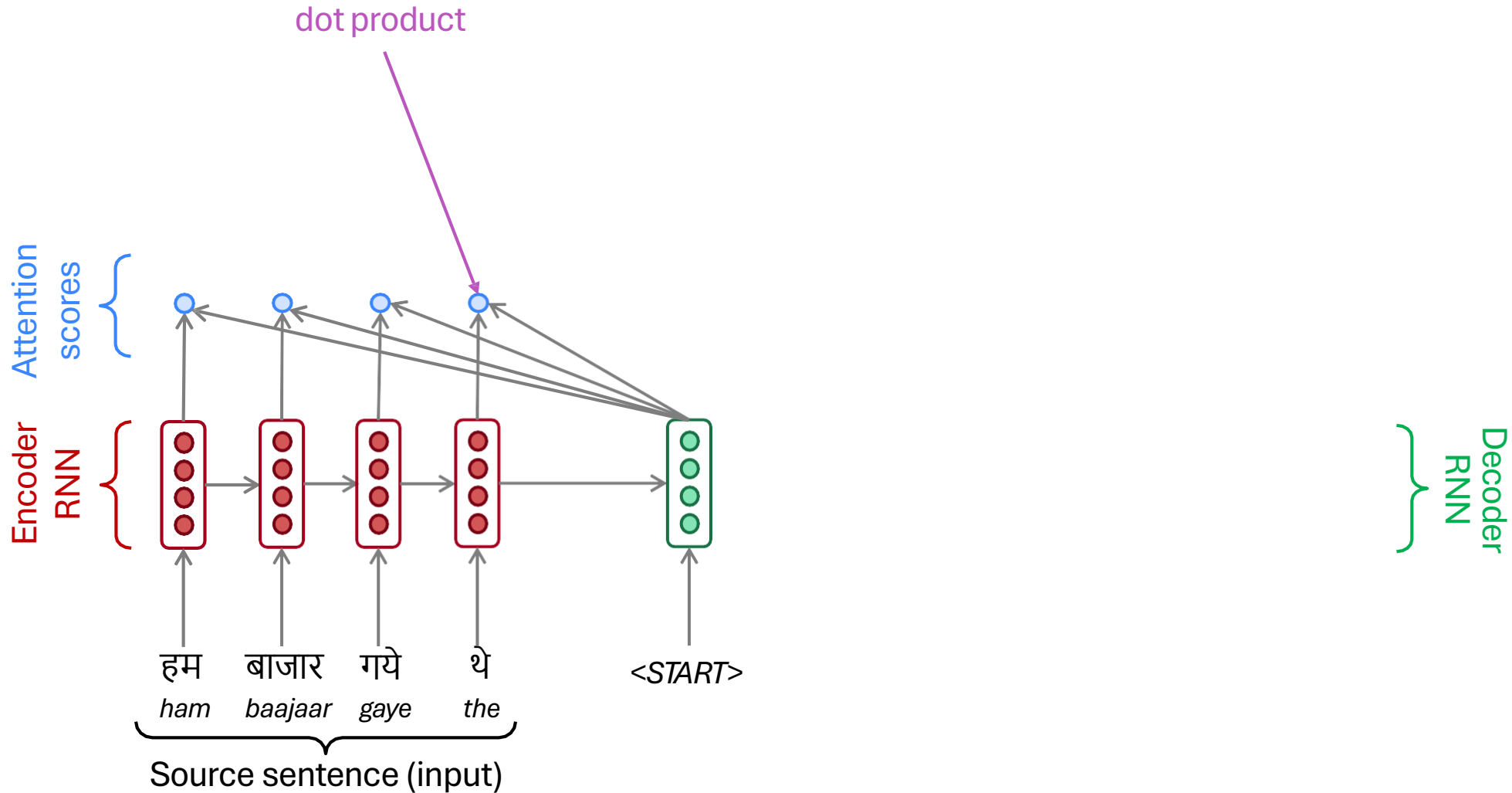हम *ham*    बाजार *baajaar*    गये *gaye*    थे *the*    <START>

Source sentence (input)

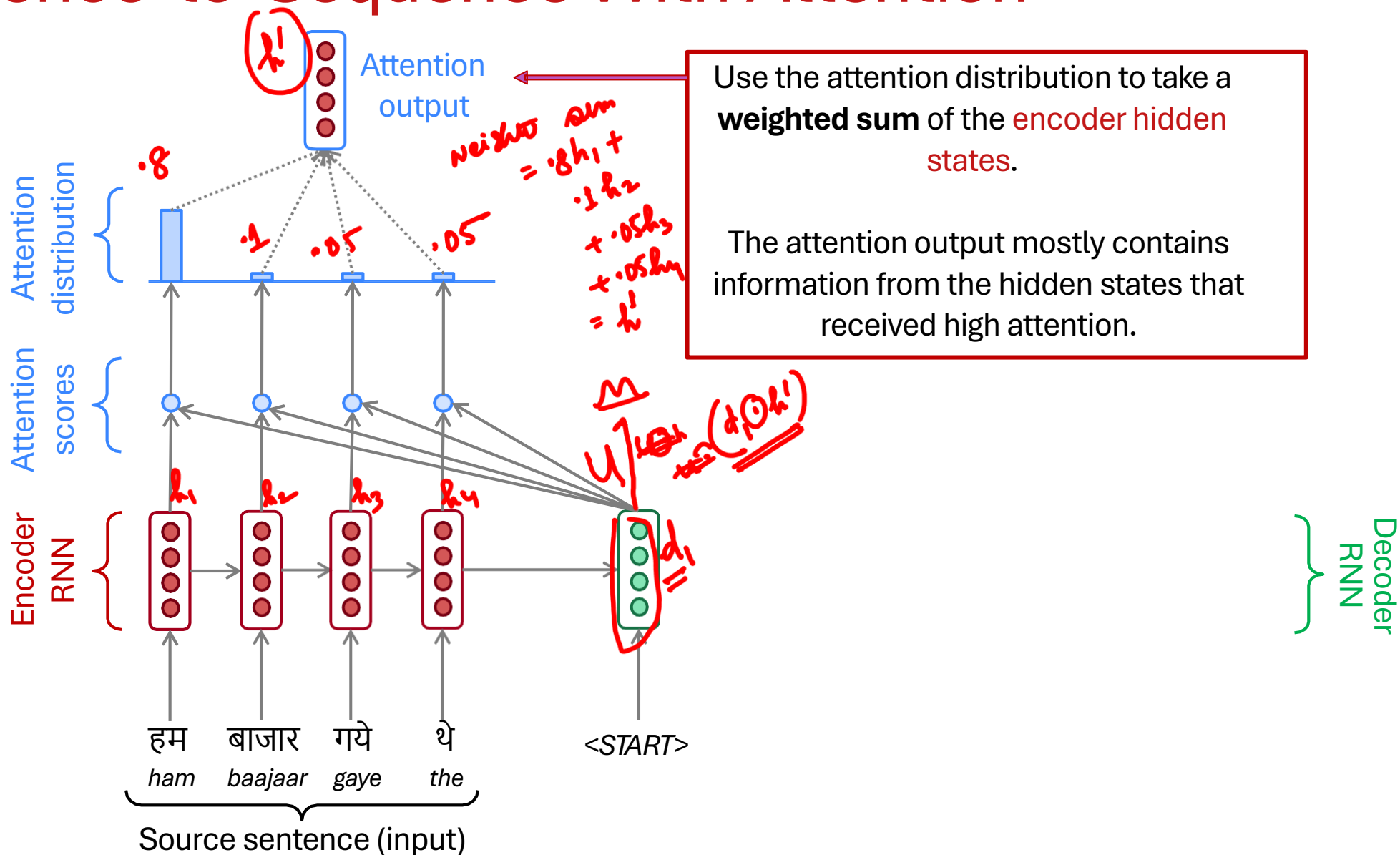# Sequence-to-Sequence With Attention

# Sequence-to-Sequence With Attention
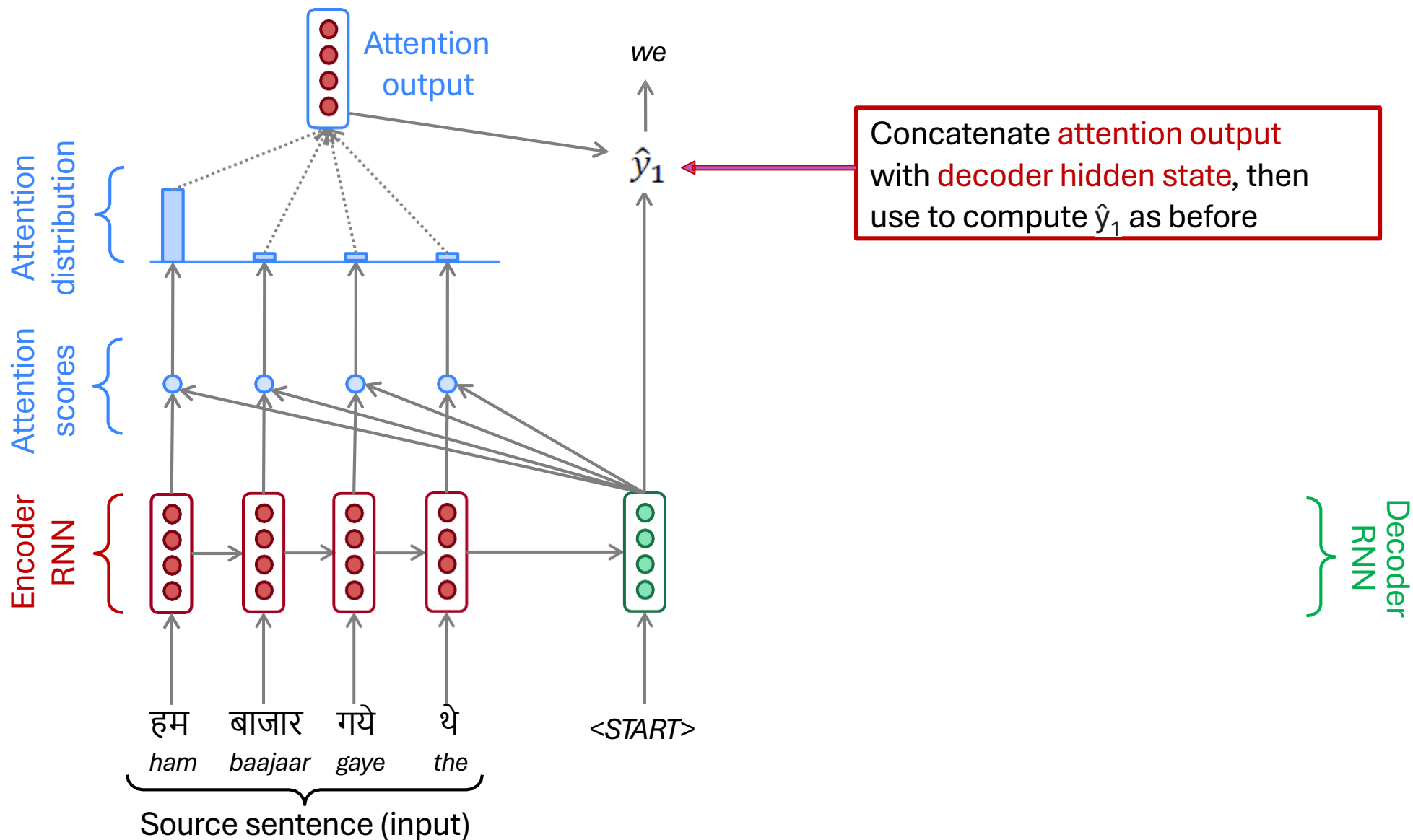
On this decoder timestep, we are mostly focusing on the first encoder hidden state

Take softmax to turn the scores into a probability distribution

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

हम
*ham*

बाजार
*baajaar*

गये
*gaye*

थे
*the*

<START>

Source sentence (input)

# Sequence-to-Sequence With Attention



Attention output

$h^1$

Use the attention distribution to take a **weighted sum** of the encoder hidden states.

The attention output mostly contains information from the hidden states that received high attention.

Attention distribution

.8    .1    .05    .05

weighted sum
$= .8h_1 + .1h_2 + .05h_3 + .05h_4 = h^1$

Attention scores

$h_1$  $h_2$  $h_3$  $h_4$

$u_i^t = f(d \odot h^i)$

$d_1 =$

Encoder RNN

Decoder RNN

हम    बाजार    गये    थे        <START>
ham   baajaar  gaye   the

Source sentence (input)

# Sequence-to-Sequence With Attention



Attention output

*we*

$\hat{y}_1$

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

Concatenate attention output with decoder hidden state, then use to compute $\hat{y}_1$ as before

हम बाजार गये थे

*ham* *baajaar* *gaye* *the*

*<START>*

Source sentence (input)

# Sequence-to-Sequence With Attention

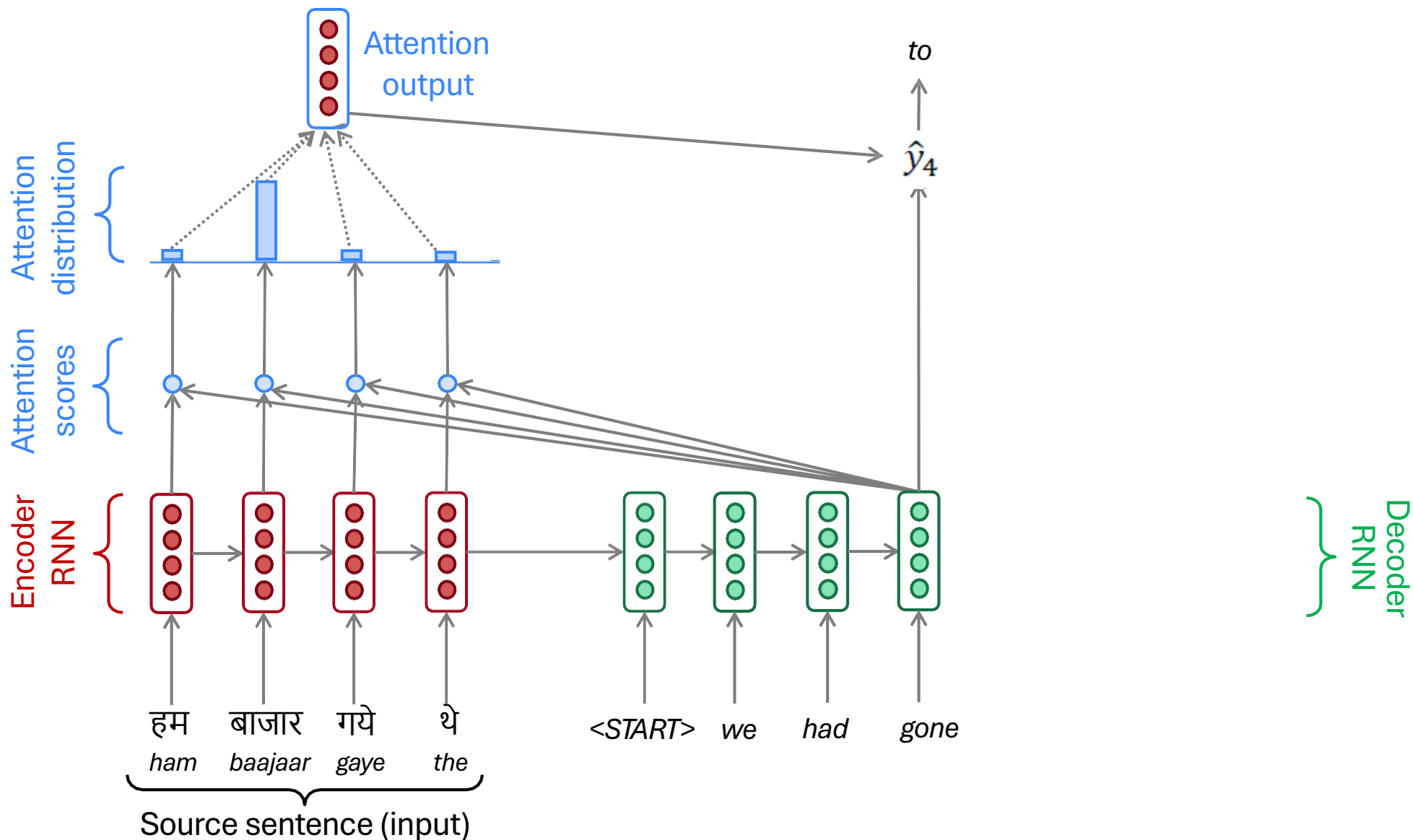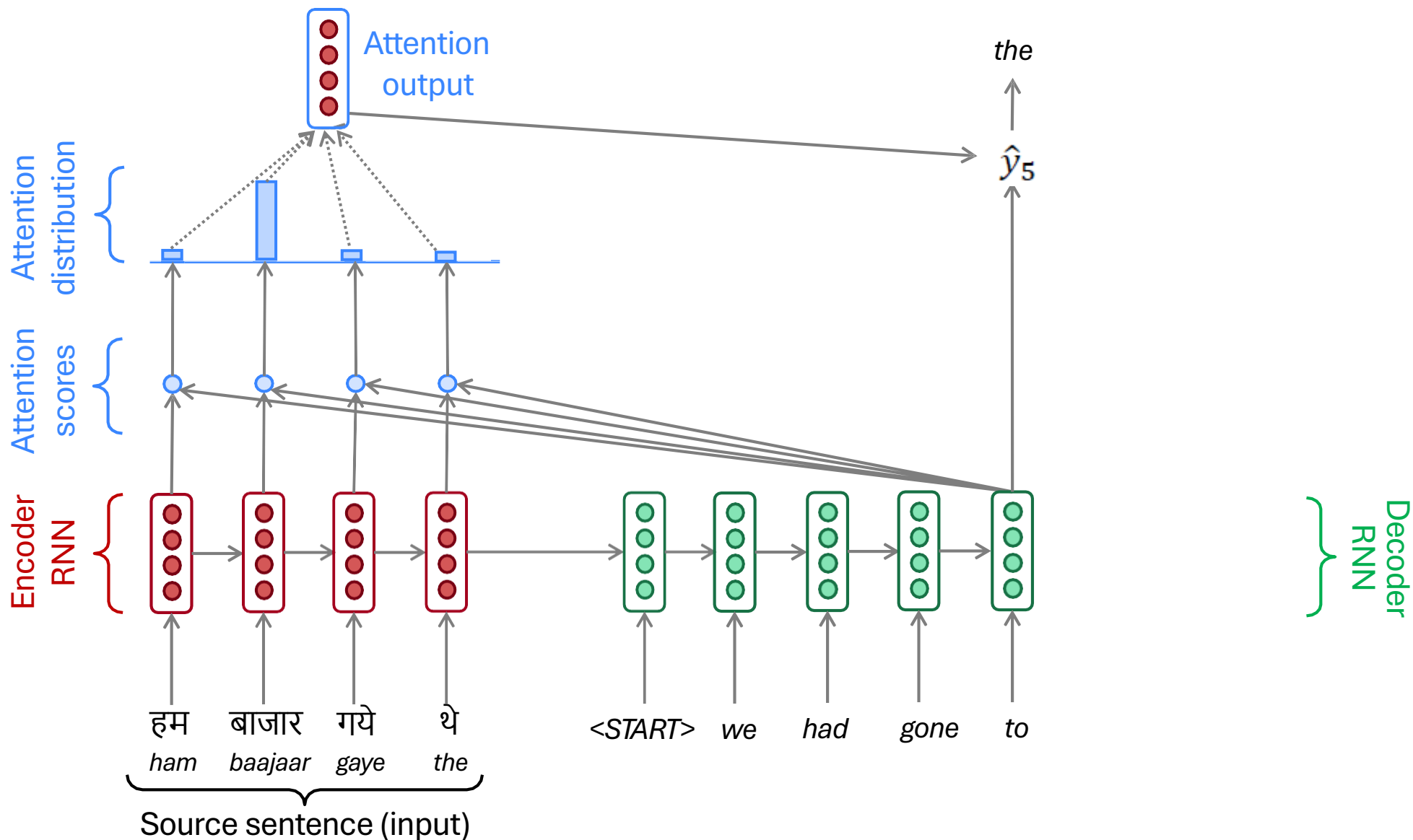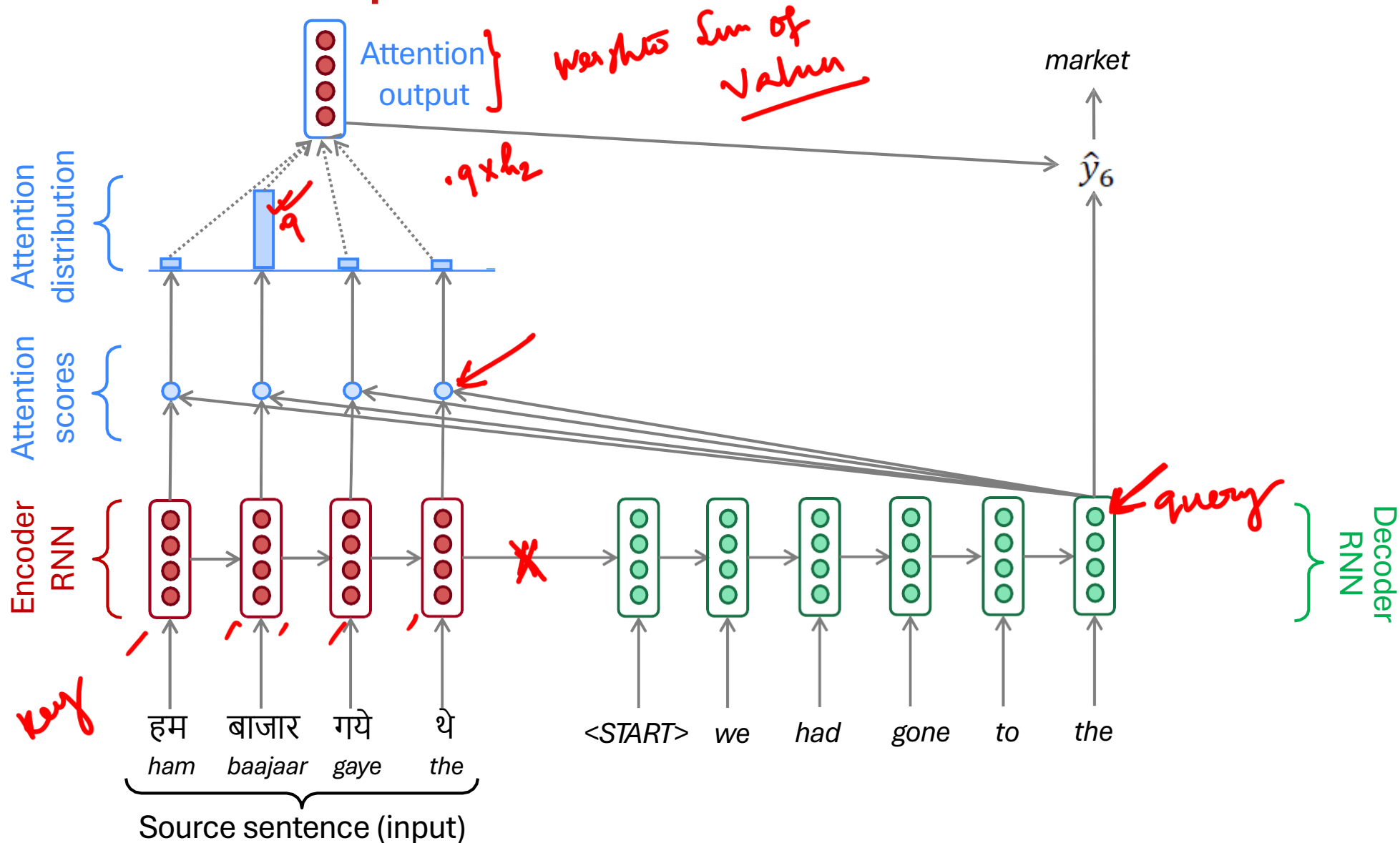# Sequence-to-Sequence With Attention

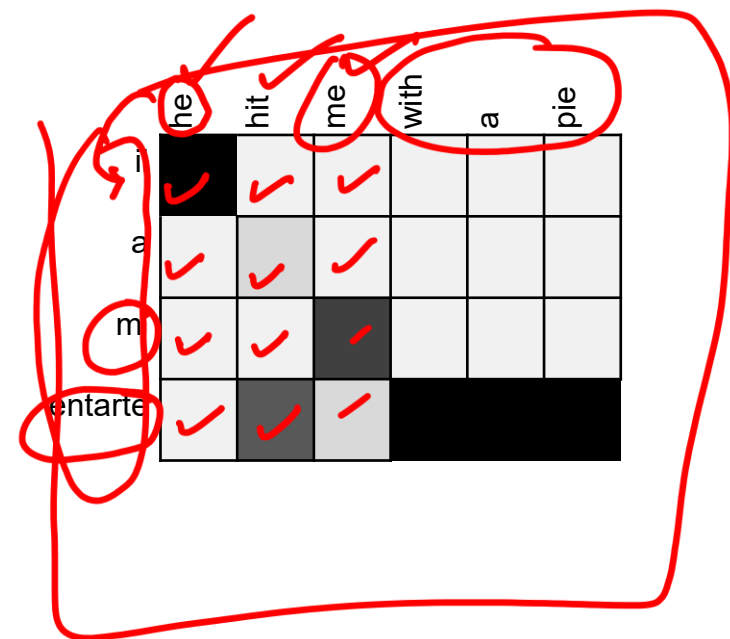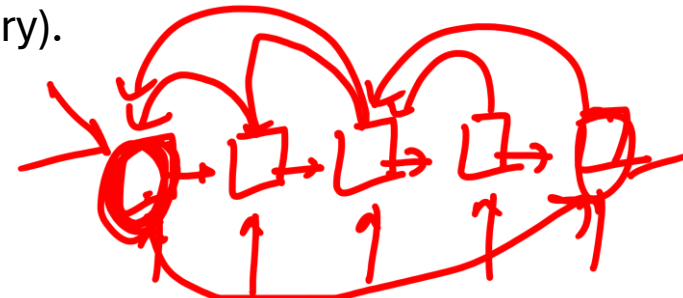# Sequence-to-Sequence With Attention



Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

$\hat{y}_4$

to

हम
ham

बाजार
baajaar

गये
gaye

थे
the

&lt;START&gt;   we   had   gone

Source sentence (input)

# Sequence-to-Sequence With Attention

# Sequence-to-Sequence With Attention



Attention output } Weighted Sum of Values

$.9 \times h_2$

Attention distribution

.9

Attention scores

Encoder RNN

$\hat{y}_6$

market

query

key

Decoder RNN

| हम | बाजार | गये | थे |
|---|---|---|---|
| *ham* | *baajaar* | *gaye* | *the* |

Source sentence (input)

<START>   we   had   gone   to   the

# Attention is Great

- Attention significantly improves NMT performance
  - It's very useful to allow decoder to focus on certain parts of the source
- Attention solves the bottleneck problem
  - Attention allows decoder to look directly at source; bypass bottleneck
- Attention helps with vanishing gradient problem
  - Provides shortcut to faraway states
- Attention provides some interpretability
  - By inspecting attention distribution, we can see what the decoder was focusing on
  - We get (soft) alignment for free!
  - This is cool because we never explicitly trained an alignment system
  - The network just learned alignment by itself

# Attention is a *General* Deep Learning Technique

- We've seen that attention is a great way to improve the sequence-to-sequence model for Machine Translation.

- <u>However</u>: You can use attention in many architectures (not just seq2seq) and many tasks (not just MT)

- More general definition of attention:

  - Given a set of vector *values*, and a vector *query*, attention is a technique to compute a weighted sum of the values, dependent on the query.

- We sometimes say that the query *attends to* the values.

- For example, in the seq2seq + attention model, each decoder hidden state (query) *attends to* all the encoder hidden states (values).

- **Intuition**:

  - The weighted sum is a *selective summary* of the information contained in the values, where the query determines which values to focus on.

  - Attention is a way to obtain a *fixed-size representation of an arbitrary set of representations* (the values), dependent on some other representation (the query).

75

# Variants of Attention

- Original formulation: $a(\mathbf{q}, \mathbf{k}) = w_2^T \tanh(W_1[\mathbf{q}; \mathbf{k}])$

- Bilinear product: $a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T W \mathbf{k}$      Luong et al., 2015

- Dot product: $a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T \mathbf{k}$      Luong et al., 2015

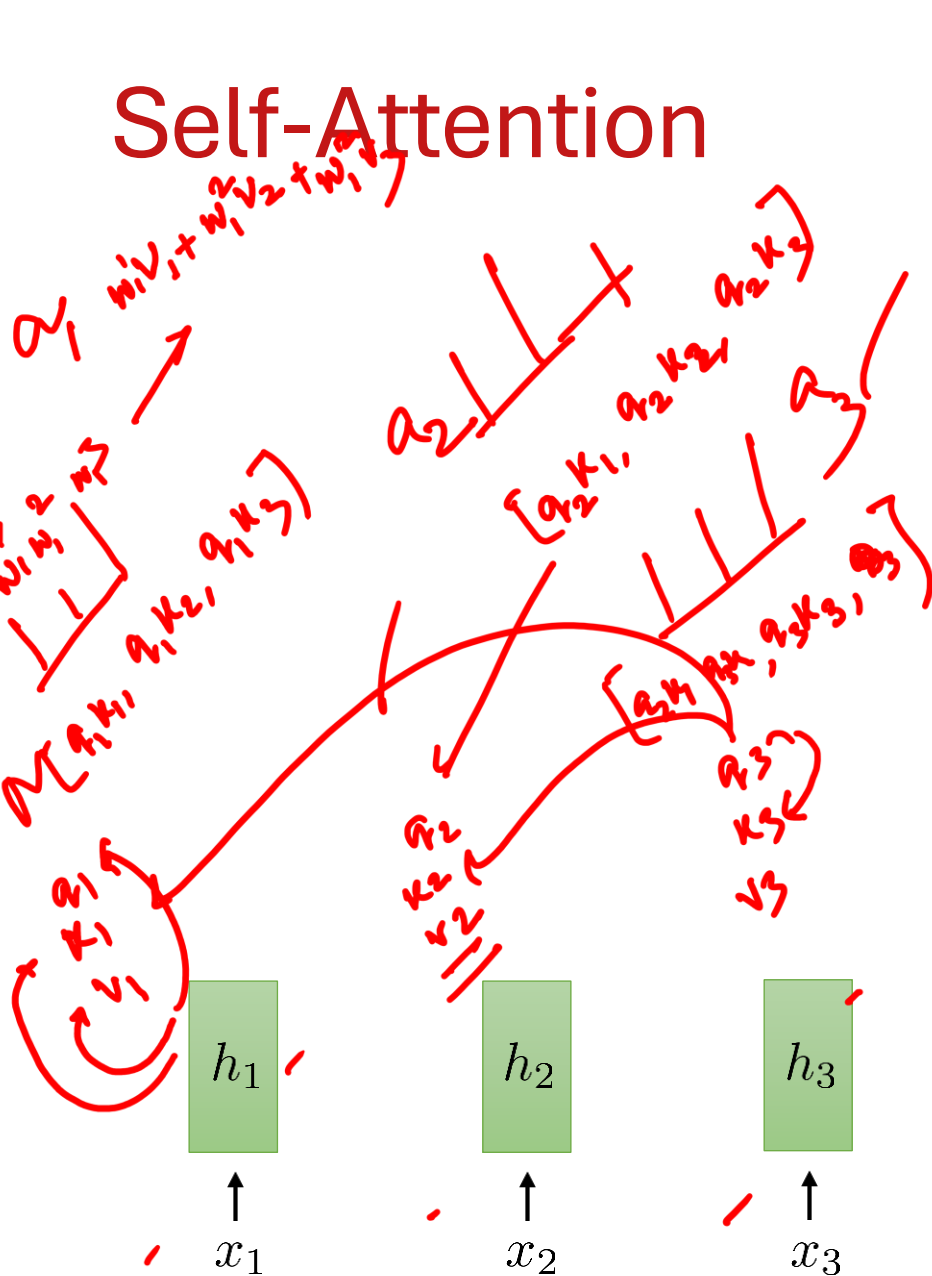- Scaled dot product: $a(\mathbf{q}, \mathbf{k}) = \dfrac{\mathbf{q}^T \mathbf{k}}{\sqrt{|\mathbf{k}|}}$      Vaswani et al., 2017

# Self-Attention



this is *not* a recurrent model!
but still weight sharing:

$$h_t = \sigma(W x_t + b)$$

shared weights at all time steps

(or any other nonlinear function)

Handwritten annotations:

$$a_1 = W_q h_1$$
$$k_1 = W_k h_1$$
$$v_1 = W_v h_1$$

$W_q$, $W_k$, $W_v$

$h_1$  $h_2$  $h_3$

$x_1$  $x_2$  $x_3$