

Explain the GRU in recNN

Below is a **clean, compact but technically precise** explanation of the **GRU (Gated Recurrent Unit)** in recurrent neural networks. I'll walk through:

1. Why GRU exists
 2. Its gates and equations
 3. Its intuition
 4. Comparison to LSTM
 5. A small numerical toy example
-

1. Why GRU Exists

Classic RNNs suffer from:

- **Vanishing gradients** → difficulty learning long-range dependencies
- **Coupling of memory and hidden state** → no mechanism to decide what to remember or forget

GRUs were introduced (Cho et al., 2014) to be a **simpler, faster alternative to LSTMs** while still solving the above problems.

GRU = RNN + **gates to control information flow**, but **fewer gates and no separate cell state** (unlike LSTM).

2. GRU Architecture & Equations

2.1 Gates

GRU has only **two gates**:

- **Update gate z_t** :
How much of the **past** to keep vs. **new** to write.
 - **Reset gate r_t** :
How much to ignore previous state when forming candidate memory.
-

2.2 Equations (standard GRU)

Let:

- $x_t \in \mathbb{R}^{d_x}$ = input at time t
- $h_{t-1} \in \mathbb{R}^{d_h}$ = previous hidden state

Update gate

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z)$$

Reset gate

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

Candidate hidden state

$$\tilde{h}_t = \tanh(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h)$$

Notice:

- If $r_t = 0$, previous state is **ignored** → good for resetting memory.
- If $r_t = 1$, it behaves like a normal RNN.

Final hidden state

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

This is the **core**:

- If $z_t \approx 0$: keep **old** memory
- If $z_t \approx 1$: use the **new** memory



3. Intuition Inside the GRU

Reset gate r_t

Should I **forget old information** when computing the new candidate?

Use case:

- Sentence: "I grew up in France ... **I speak French**"
When predicting "French", previous "I grew up in France" should matter.
- Sentence: "Today is cold ... **I ate a burger**"
Past climate should not matter → reset gate low.

Update gate z_t

Should I **update the memory** or **keep the existing memory**?

Use case:

- When the network should carry long-term context
→ Large z_t
- When new input overwrites memory
→ Small z_t

4. GRU vs LSTM

Feature	LSTM	GRU
Gates	3 (input, forget, output)	2 (update, reset)
Cell state?	Yes (separate)	No (hidden = memory)
Parameters	More	~25% fewer
Speed	Slower	Faster
Long dependencies	Excellent	Very good
Empirical performance	Often best	Comparable, simpler

5. Small Toy Example

Let's do a **1-dimensional GRU** to see updates explicitly.

Assume:

- $h_{t-1} = 0.6$
- input $x_t = 1.0$

Let gates be:

- $z_t = \sigma(1.2) \approx 0.768$
- $r_t = \sigma(-1.0) \approx 0.268$

Candidate:

$$\tilde{h}_t = \tanh(W_h x_t + U_h(r_t h_{t-1}))$$

Say:

- $W_h x_t = 0.5$

- $U_h(r_t h_{t-1}) = 0.1$

Then:

$$\tilde{h}_t = \tanh(0.6) \approx 0.537$$

Final state:

$$\begin{aligned} h_t &= (1 - z_t)h_{t-1} + z_t\tilde{h}_t \\ &= (1 - 0.768)(0.6) + 0.768(0.537) \\ &= 0.232 \cdot 0.6 + 0.768 \cdot 0.537 \\ &= 0.139 + 0.412 \approx 0.551 \end{aligned}$$

So the new hidden state is **0.551**.

Because the update gate was large (0.768), the network mostly uses the **new candidate**.

✓ Summary

- GRU is a gated RNN variant designed to handle long-term dependencies.
- It uses **two gates**: update (z) and reset (r).
- It **mixes old and new information smoothly**, producing stable gradients.
- It is **simpler and faster than LSTM**, with comparable performance.

If you want, I can also produce:



GRU vs LSTM full derivation



PyTorch GRU code examples (raw or nn.GRU)



Diagram of the GRU cell



PDF version formatted as lecture notes

Just tell me!

