# Unraveling LSTMs: A Journey Through Questioning

## Introduction

When learning about deep learning architectures, it's common to focus on implementation details without fully grasping what happens under the hood. This blog post captures an insightful journey of questioning and discovery—where we aimed to uncover the nuances of how **stacked LSTMs** behave, particularly when dealing with multiple layers and `return_sequences=False`.

Through a series of back-and-forth inquiries, we explored how information flows in LSTMs, how they differ from Dense layers, and when LSTMs are truly leveraging their memory mechanisms.

---

## 1. Understanding LSTM Cells, Units, and Neurons

One of the first questions that arose was about terminology: **What exactly do we mean by LSTM cells, units, and neurons?**

- **LSTM cell**: A single computational unit responsible for processing a single timestep of an input sequence.

- **LSTM unit**: Equivalent to the number of hidden dimensions in an LSTM layer (e.g., `LSTM(50)` means 50 units, each with its own cell state and hidden state).

- **Neurons**: A more general term used in traditional neural networks; in LSTMs, a "unit" functions similarly to a neuron but operates with additional gating mechanisms.

Key takeaway: **Each LSTM cell processes one timestep at a time, passing information through time using hidden and cell states.**

---

## 2. How Does Information Flow in Stacked LSTMs?

### Scenario: `return_sequences=False` in the First LSTM Layer

A major point of exploration was understanding what happens when `return_sequences=False` in the first LSTM layer of a stacked model. Let's visualize this with an example.

```python
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dense

model = Sequential([
    LSTM(50, return_sequences=False, input_shape=(10, 5)),  # First LSTM layer
    LSTM(10, return_sequences=False),  # Second LSTM layer
    Dense(1, activation='sigmoid')  # Output layer
])

model.summary()
```

Here's what's happening at each layer:

1. **First LSTM Layer (50 units, `return_sequences=False`)**

   - Input shape: `(batch_size, 10, 5)` (10 timesteps, 5 features per timestep)

   - **Processes the entire sequence** but **outputs only the final hidden state**.

   - Output shape: `(batch_size, 50)` → **No more time steps left!**

2. **Second LSTM Layer (10 units, `return_sequences=False`)**

   - Input shape: `(batch_size, 50)` (single vector, no timesteps left).

   - **This means it is NOT processing a sequence anymore**—it behaves more like a Dense layer.

   - Output shape: `(batch_size, 10)`.

3. **Dense Layer**

   - Input shape: `(batch_size, 10)`.

   - Output shape: `(batch_size, 1)` (final prediction).

## Key Realization

- **Since the first LSTM layer collapses the sequence, the second LSTM layer no longer processes time steps—it behaves like a fully connected layer with LSTM-style computations.**

- This means **LSTM's memory capabilities are not being used in the second layer the way they are in the first layer**.

# 3. Does the Second LSTM Layer Act Like a Dense Layer?

This led to the next crucial question: **If the second LSTM layer no longer processes a sequence, does it behave just like a Dense layer?**

The answer: **Not exactly, but similar.**

| Layer Type | Operations | Memory Retention? | Output Shape |
|---|---|---|---|
| **LSTM(10, return_sequences=False)** | Uses forget, input, and output gates | ✅ Yes (but minimal effect here) | `(batch_size, 10)` |
| **Dense(10)** | Simple `Wx + b` transformation | ❌ No memory | `(batch_size, 10)` |

- While both layers take a **fixed-length vector as input** (instead of a sequence), an LSTM still applies its **gating mechanisms**.

- **However, if the LSTM is not effectively leveraging memory, its behavior becomes close to a Dense layer.**

## Final Experiment: Replace LSTM with Dense

To validate this, we can swap the second LSTM layer with a Dense layer:

```
model = Sequential([
    LSTM(50, return_sequences=False, input_shape=(10, 5)),
    Dense(10, activation='relu'),
    Dense(1, activation='sigmoid')
])
```

Would this perform just as well? **If memory is not important for the second layer, then yes!** If performance drops, it suggests that the second LSTM layer was still

utilizing gating mechanisms to retain some information.

# 4. Conclusion: What Did We Uncover?

Through these back-and-forth questions, we uncovered:

✅ `return_sequences=False` **collapses the sequence, passing only the final hidden state forward.**

✅

**A second LSTM layer without sequences behaves more like a Dense layer but still applies LSTM gating.**

✅

**If memory is not important at this stage, replacing the second LSTM with a Dense layer may give similar results.**

✅

**Stacking LSTMs only makes sense when** `return_sequences=True` **is used, allowing sequential information to propagate across layers.**

## Final Thought 💡

When designing LSTM-based models, always ask:

- **Is my second LSTM layer actually processing a sequence, or is it acting like a Dense layer?**

- **Do I need memory at this stage, or should I replace it with a simple Dense layer?**

- **Is** `return_sequences=True` **set where necessary to allow sequential processing?**

Understanding these nuances will help you **build more efficient models without unnecessary complexity**! 🚀

## What's Next?

Now that we've explored stacked LSTMs, the next step is to experiment:
🔷 Try different values of
`return_sequences`
🔷 Replace LSTMs with Dense layers and compare performance
🔷 Investigate bidirectional LSTMs for better sequence learning

Have more questions? Keep questioning and uncovering! 🔍