# TTT - Learning to (Learn at Test Time)



Figure 4. **Top**: A generic sequence modeling layer expressed as a hidden state that transitions according to an update rule. All sequence modeling layers can be viewed as different instantiations of three components in this figure: the initial state, update rule and output rule. **Bottom**: Examples of sequence modeling layers and their instantiations of the three components. The naive TTT layer was shown in Figure 1. Self-attention has a hidden state growing with context, therefore growing cost per token. Both the naive RNN and TTT layer compress the growing context into a hidden state of fixed size, therefore their cost per token stays constant.

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□                                      TTT□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□                          compression heuristic
□□□□□□□□□□□□□□□□token□□□□□□□□□□□□□□□□□□□□
Transformer□KV cache□□□□□□□□□□□□□          Manba
□□□□□□□□□□□□□□□□□□□□□□

## □□□□

□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□

████████████████████████████ f████ ██ token███████████████ f████████
████████████ ℓ███████████

████████ f██████████████████████████

████████████████████████████

████████████ -█████ TTT██

## ██

████ TTT████████████████████ Transformer
█████████████████ token███████████████

██████████████████ outer loop███████ TTT███ W████ inner loop
████████████████ W██ f
████████████████

## ██

```python
class TTT_Layer(nn.Module):
  def __init__(self):
    self.task = Task()

  def forward(self, in_seq):
    state = Learner(self.task)
    out_seq = []
    for tok in in_seq:
      state.train(tok)
      out_seq.append(state.predict(tok))
    return out_seq

class Task(nn.Module):
  def __init__(self):
    self.theta_K = nn.Param((d1, d2))
    self.theta_V = nn.Param((d1, d2))
    self.theta_Q = nn.Param((d1, d2))

  def loss(self, f, x):
    train_view = self.theta_K @ x
    label_view = self.theta_V @ x
    return MSE(f(train_view), label_view)
```

```python
class Learner():
  def __init__(self, task):
    self.task = task
    # Linear here, but can be any model
    self.model = Linear()
    # online GD here for simplicity
    self.optim = OGD()

  def train(self, x):
    # grad function wrt first arg
    # of loss, which is self.model
    grad_fn = grad(self.task.loss)
    # calculate inner-loop grad
    grad_in = grad_fn(self.model, x)

    # starting from current params,
    # step in direction of grad_in,
    self.optim.step(self.model, grad_in)

  def predict(self, x):
    test_view = self.task.theta_Q @ x
    return self.model(test_view)
```

Figure 6. Naive implementation of a TTT layer with a linear model and online GD in the style of PyTorch. TTT_Layer can be dropped into a larger network like other sequence modeling layers. Training the network will optimize the parameters of Task in TTT_Layer, because both are subclasses of nn.Module. Since Learner is not a subclass of nn.Module, state.model is updated manually in the inner loop for each call of state.train. For simplicity, we sometimes overload model as model.parameters.

□□□□□□□□□□□□□□□□□ **Task**□□□□□□□□□□□□□□ **theta_K theta_V**□□□□□□ **sequence**□□□□ **theta_K theta_V theta_Q** □□□□□□□□□□□□□□

## □□

□□□□ Pile□□□□ 2k□ 8k□□□□□□□□□□□□□ Pile□□□□□□□□□ LLM□□□□□□□□□

TTT-MLP□ M□□□□□ FLOP□□□□□□□□□□ TTT-MLP□□□□□□□□□ TTT-Linear □□□□□□□□□ FLOP□□□□□□□□□□□□

□ 8k□□□□□ TTT-Linear□ M□□ TTT-MLP□ M□□□□□□□□ Mamba□□□□□ Transformer □□□ TTT-MLP□ T□□□□□□ Mamba□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ TTT□□□□ Mamba □□□□□□□

---