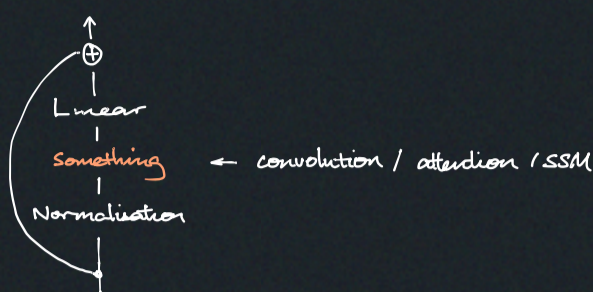# Lec 16 | Long Seq Modelling

Deep sequence models very useful

Sequence model: seq $\xrightarrow{f \circ (\cdot)}$ seq

RNN, CNN, Neural ODEs, Attention

```
    ↑
    ⊕
    |
  Linear
    |
  Something    ← convolution / attention / SSM
    |
  Normalisation
    |
```

# Baselines

RNN — recurrence cells
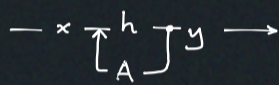+ inherently causal
- slow training
- vanishing gradient

Attention
+ parallelisable
- quadratic time training
  ↳ finite context window

→ Selective State Spaces
+ parallel, fast, linear
+ long context
+ performance

# SSM

```
— x ⟶ [ h ] ⟶ y ⟶
        [ A ]
```

signal processing
continuous, diffEQ
like RNN with continuous time,
            big hidden state,
            no activation,
            linear

$h'(t) = A h(t) + B x(t)$
         ↑            ↑ transform input
    state transition

$y(t) = C h(t) + D x(t)$
         ↑            ↑ skip connection
    project back to 1D

Biased for continuous data, less good for text

▷ Discretised

$h = \bar{A} h + \bar{B} x$
$y = \bar{C} h + \bar{D} x$

* Issues
  - no parallelisation along seq length
  - large hidden (expensive)
  - harder to train, given known future

▷ Convolution version of SSM

Equivalent to convolution
$y(t) = x(t) * K(t)$
         ↑ convolution kernal defined by A, B, C, D

+ Fast fourrier transform, near linear time

# SSSM

Linear Time Invariant — params invariant through time

SSSM — State Space Seq Model
SSSSM — Structured State Space Seq Model

Things viewed as SSM:

- RNN → state is one fixed-sized vec
                    efficient, maybe too strong
                    of compression on history
- Attention → state is cache of entire history
                    attend to all past key & value
                    good history, bad performance

# Better model

Compress, selectively remember relevant information

Doing the selection: parametrise update func on current input

Make efficient: need to tailor to GPU hardware