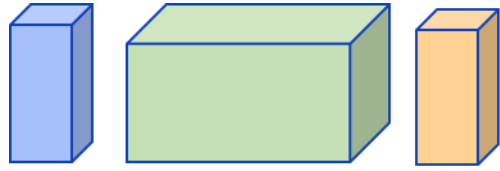


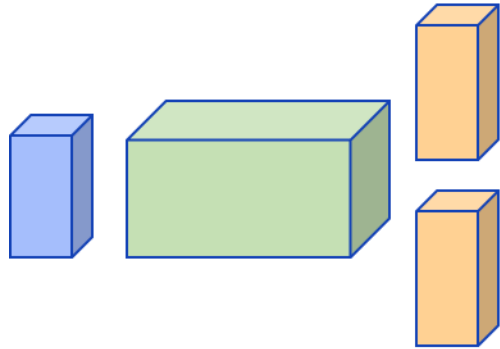
SEQUENCES AND LANGUAGE MODELS

Keith Butler

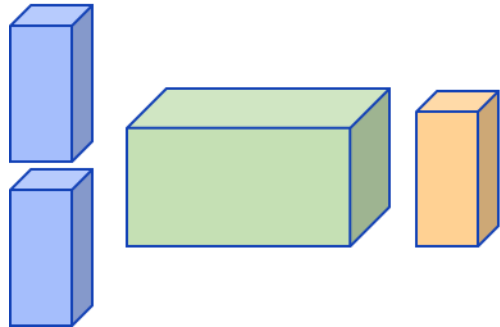
SEQUENCE DATA



One to one



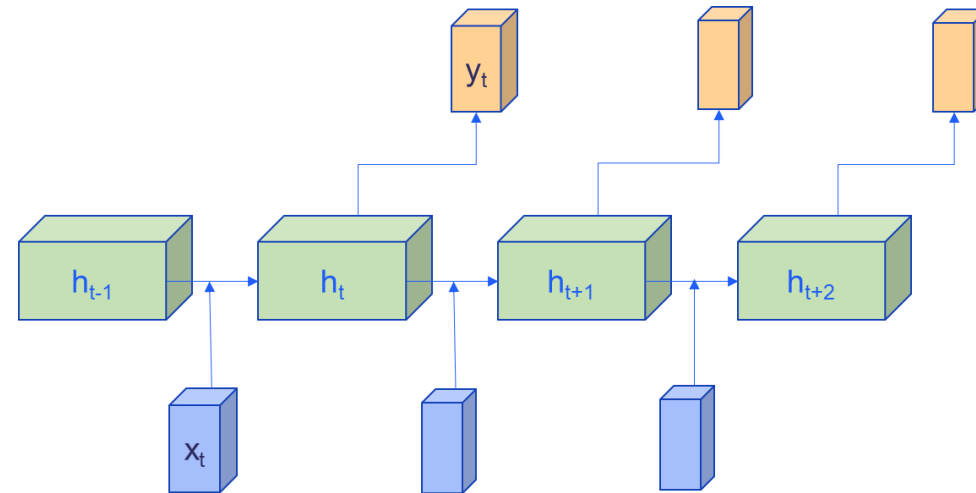
One to many



Many to one

MLPs and CNNs struggle with sequence data. There is no fixed length of input, the data can be many to one.

RECURRENT NEURAL NETWORKS (RNNs)



$$h_t = f_w(W_{hh}h_{t-1}, W_{xh}x_t)$$

ISSUES WITH RNNS

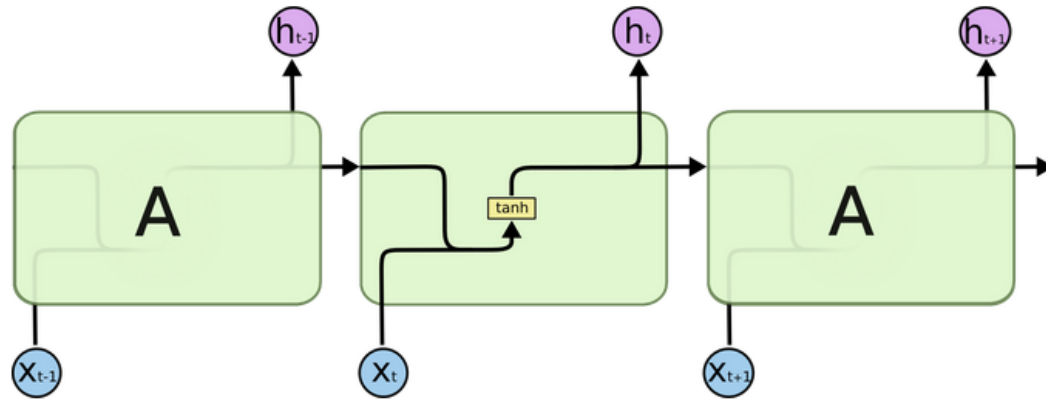
RNNs have very **short term memory** – they lose **context** quickly

The clouds are in the _____ . ✓

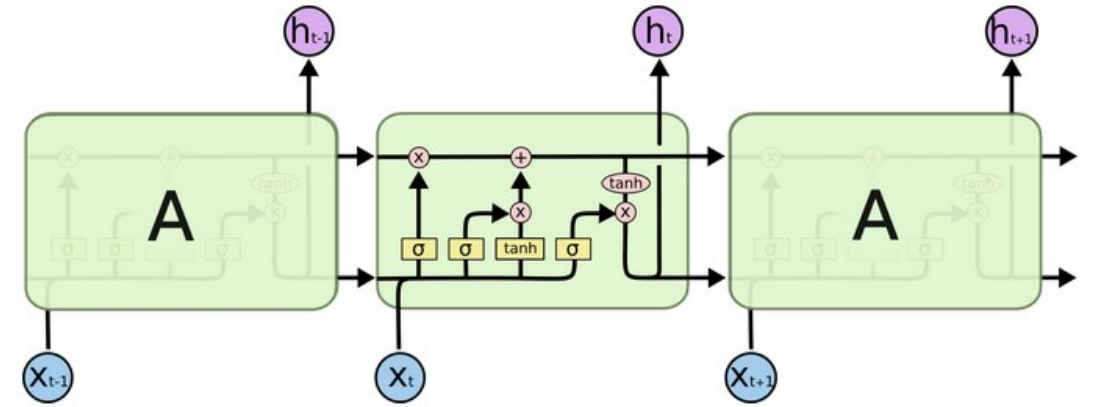
I was born in France. At the age of 16 I moved country. I have lived here since I was 21. Nonetheless, I still speak fluent _____ . ✗

INTRODUCING MEMORY

RNN



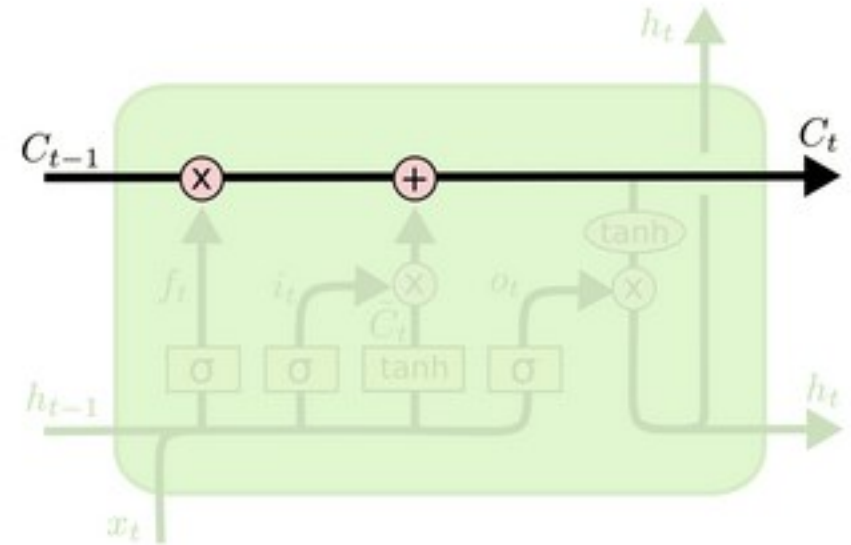
LSTM



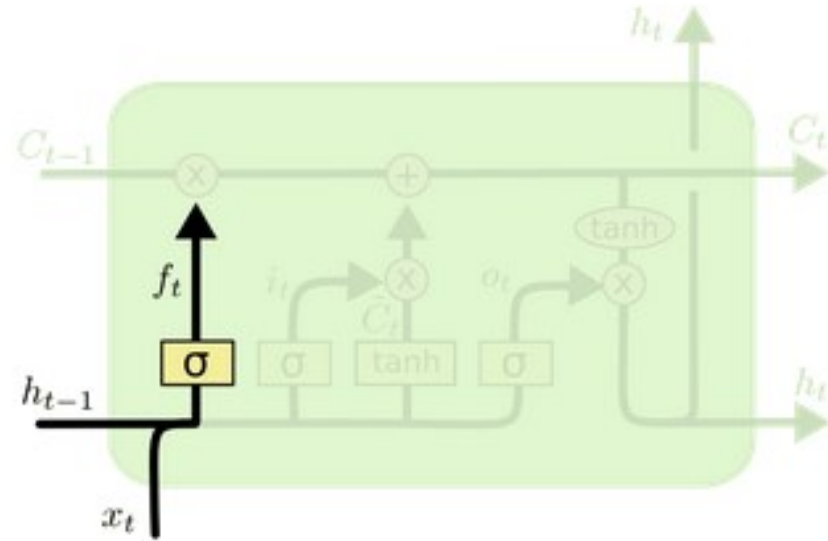
Long short term memory (LSTM) networks introduce extra memory features compared to a standard RNN

LSTM – THE MEMORY STATE

A **single channel** that runs all the way along the **sequence structure**
Only has some **minor interactions** with the rest of the network

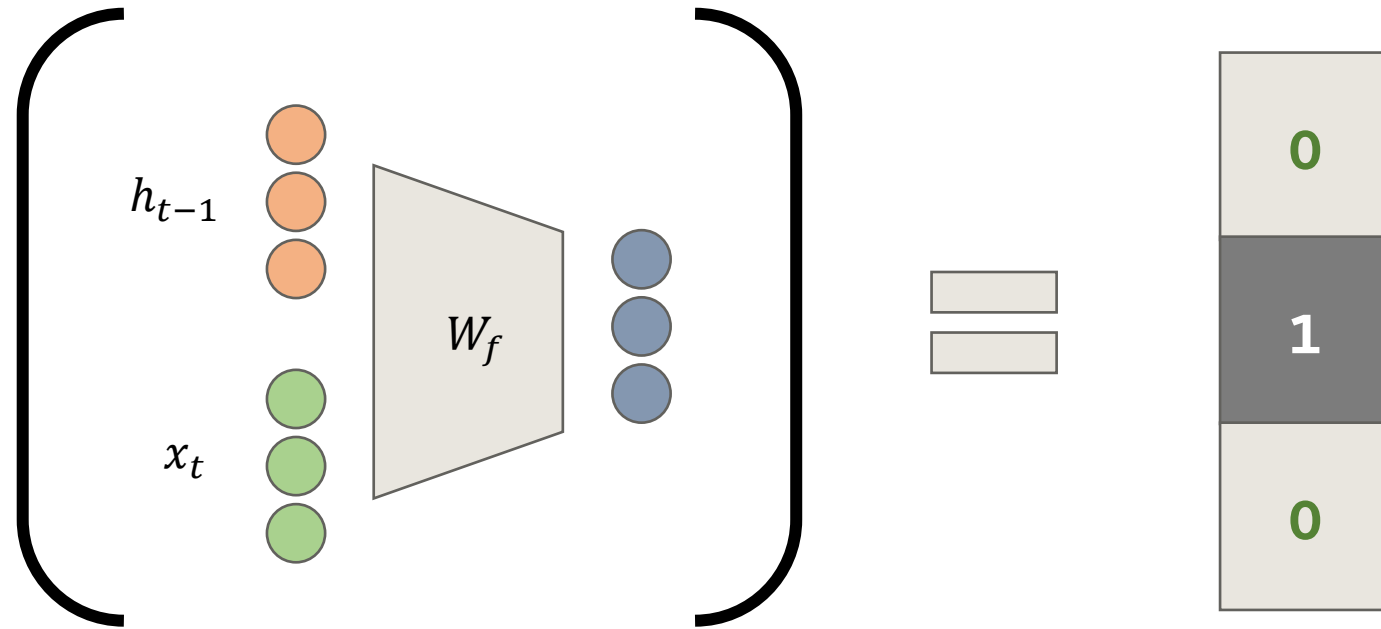
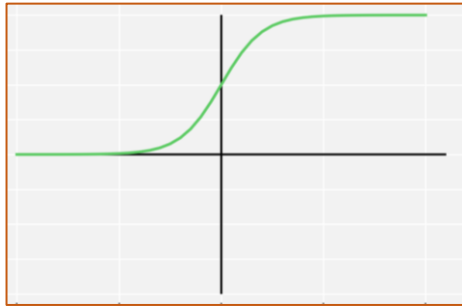


LSTM - FORGETTING



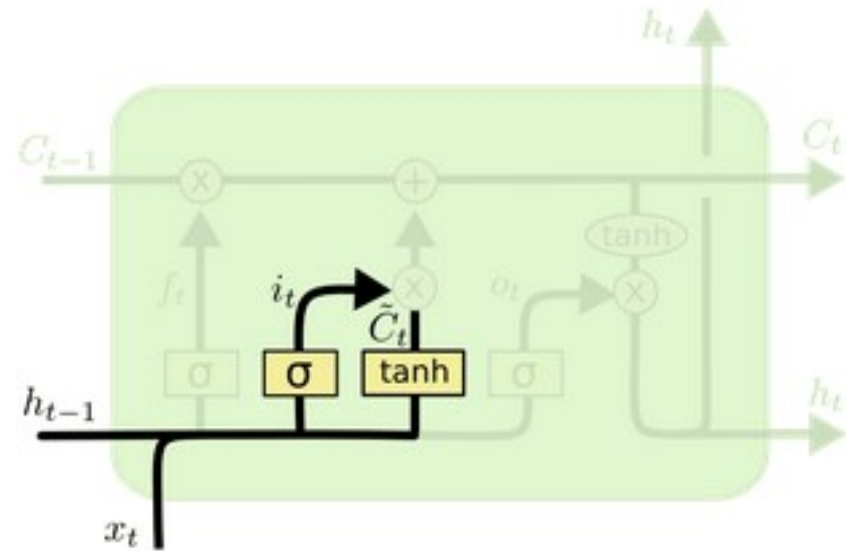
$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

FORGETTING



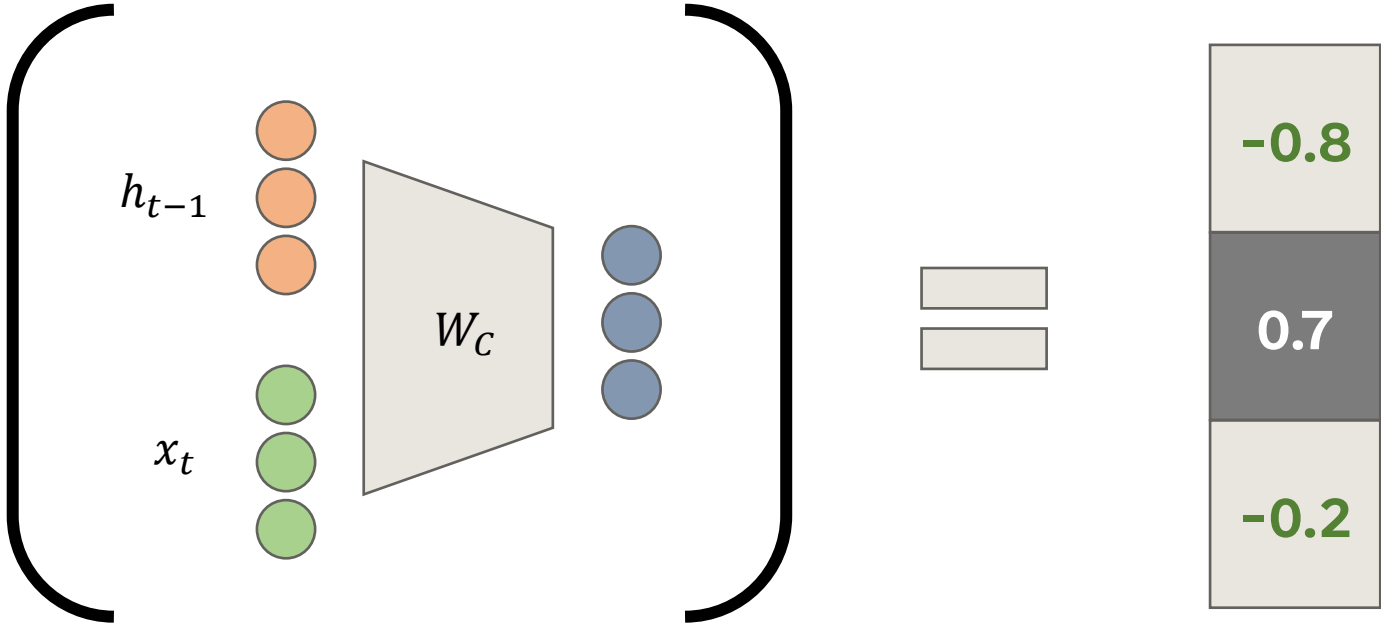
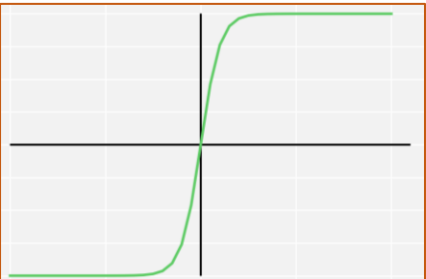
$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

LSTMS - ADDING



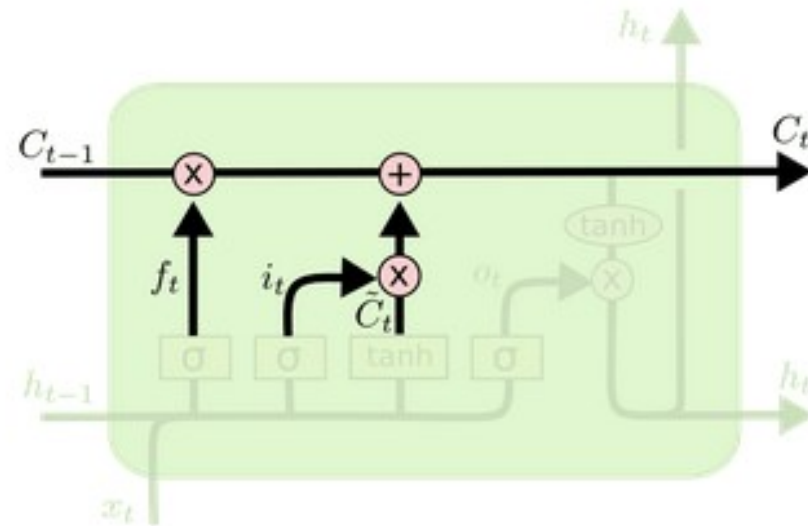
$$i_t = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$
$$\tilde{C}_t = \tanh (W_C \cdot [h_{t-1}, x_t] + b_C)$$

ADDING



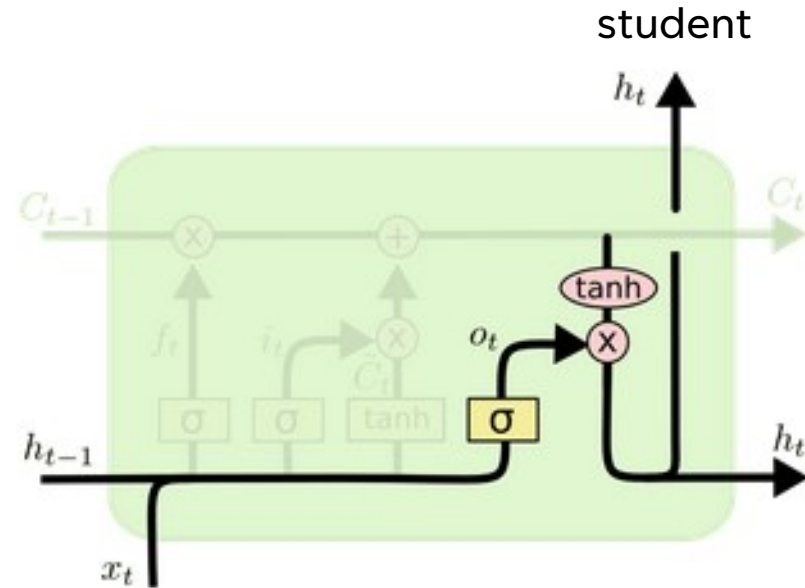
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

LSTM – UPDATING THE MEMORY



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

LSTM – GENERATE OUTPUT



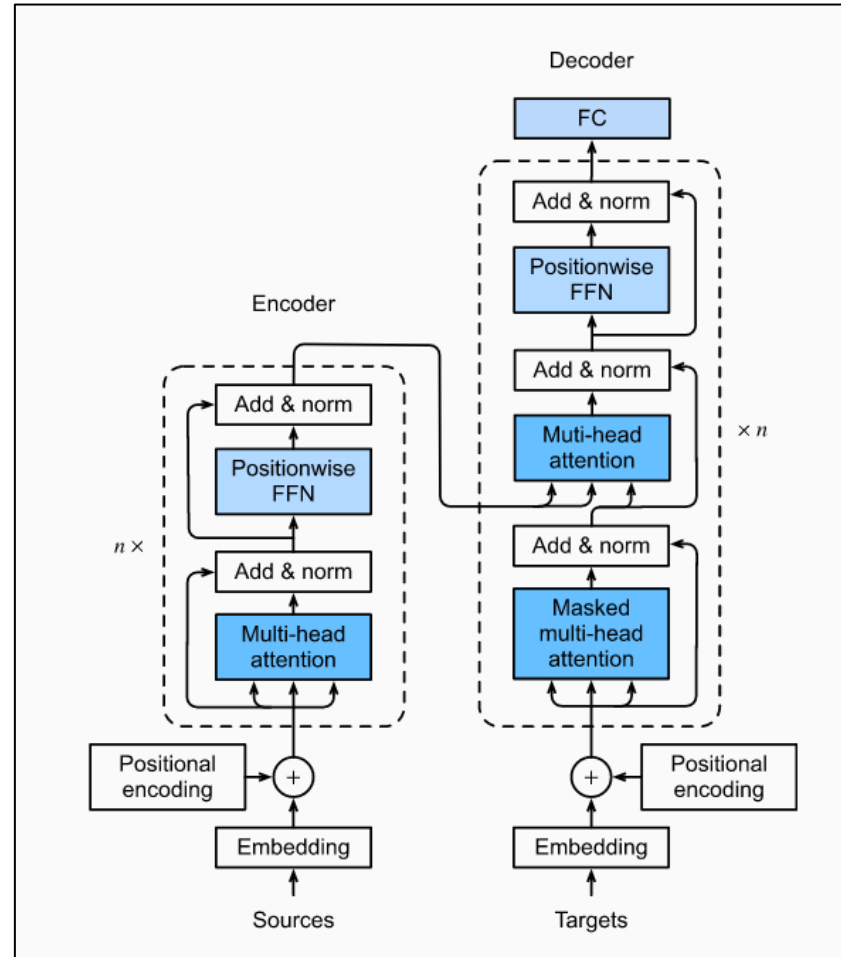
Je suis **etudiant**

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

TRANSFORMERS

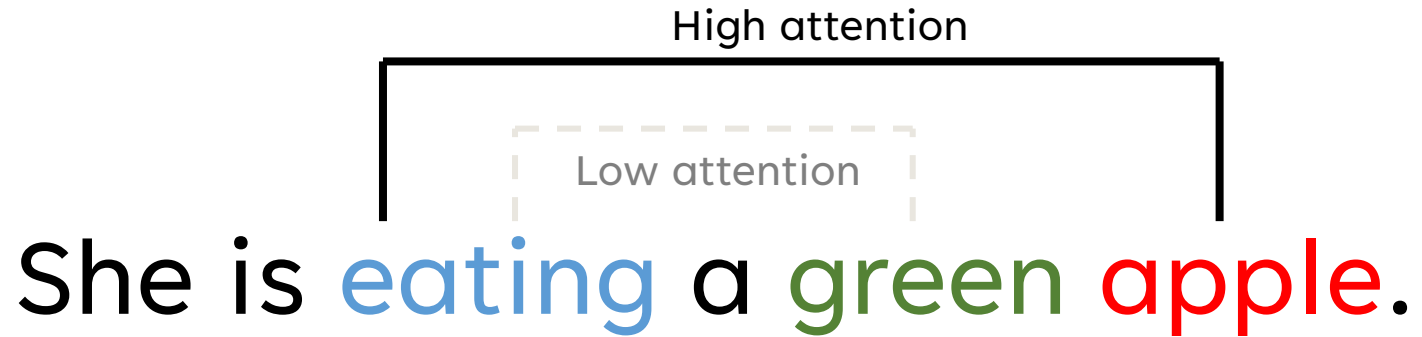
student



Je suis etudiant

I am a

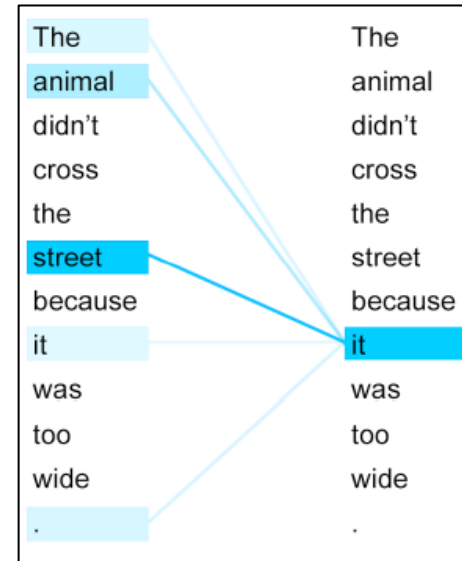
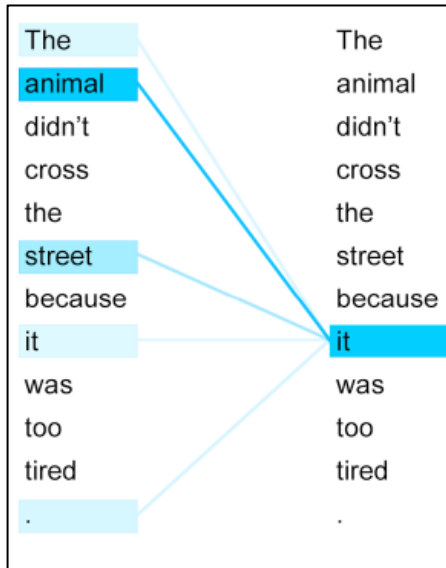
THE ATTENTION MECHANISM



We use attention to “focus” on some part of interest in an input

SELF-ATTENTION

With self-attention, each token t_n can “attend to” all other tokens of the same sequence when computing this token’s embedding x_n



HOW SELF-ATTENTION WORKS

Attention Is All You Need

$$\textit{Attention}(Q, K, V) = \textit{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$

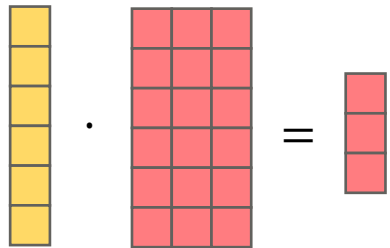
Based of the concept of **query**, **key** and **value** vectors

ATTENTION – THE INPUTS

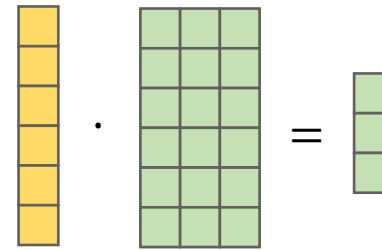
Each token has a **d-dimensional representation**

Each token also has a query and key vector q-dimensional; $q \ll d$

W_K is a matrix of learnable weights



$$\vec{x}_i \cdot W_Q = \vec{Q}_i$$



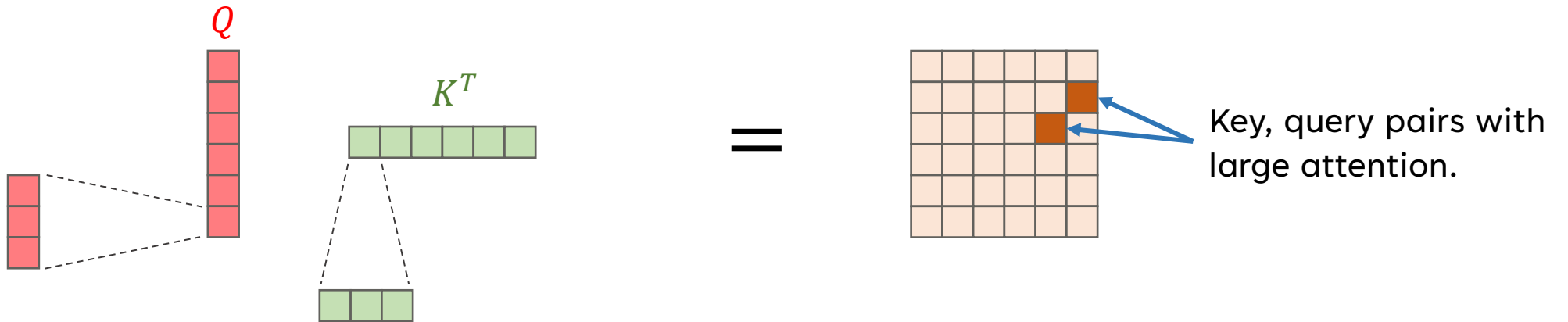
$$\vec{x}_i \cdot W_K = \vec{K}_i$$

In bed an old woman lies.

QUERY KEY MULTIPLICATION

$$\text{softmax} \left(\frac{QK^T}{\sqrt{d}} \right)$$

In bed an old woman lies.

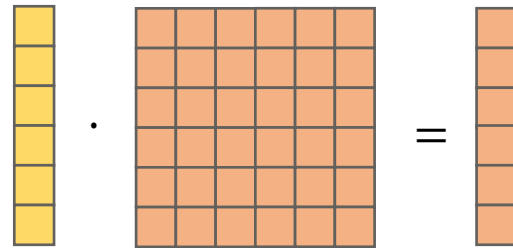


Softmax normalises the columns; root of d makes it numerically stable

NOTE –in this case each cube in **Q** and **K** has 3 dimensions

THE VALUE MATRIX

Tells you how a given token modifies another token
The resultant \vec{V} gets added to the other vector
The extent of the addition is scaled by the QK^T product



In ---- an old woman lies.

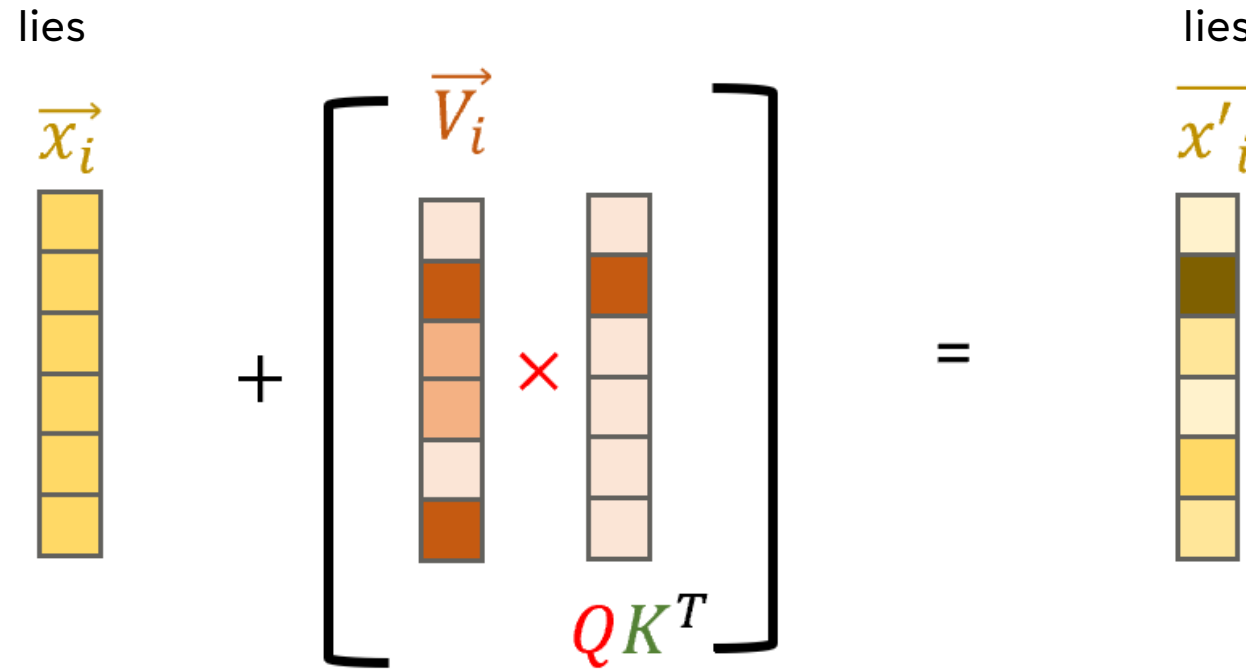
bed

court



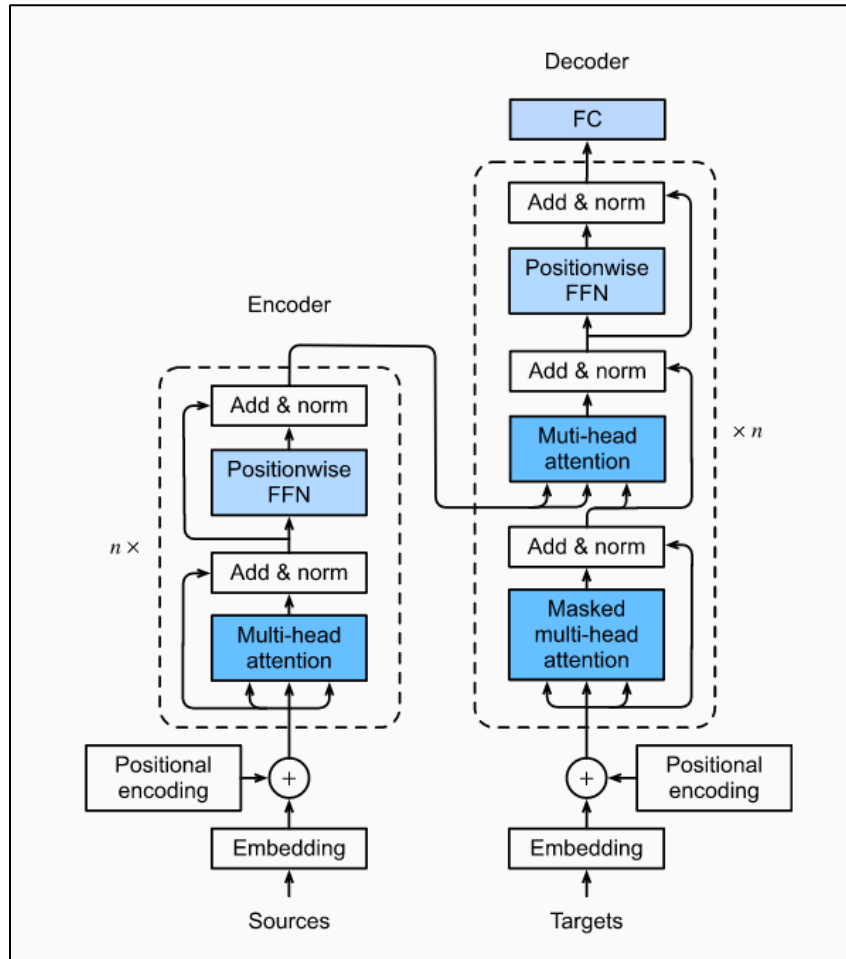
$$\vec{x}_i \cdot W_V = \vec{V}_i$$

ADDING THE VALUE TO THE EMBEDDING



The **value** is modified by the **attention** from the **QK** pair and added to the initial **embedding**

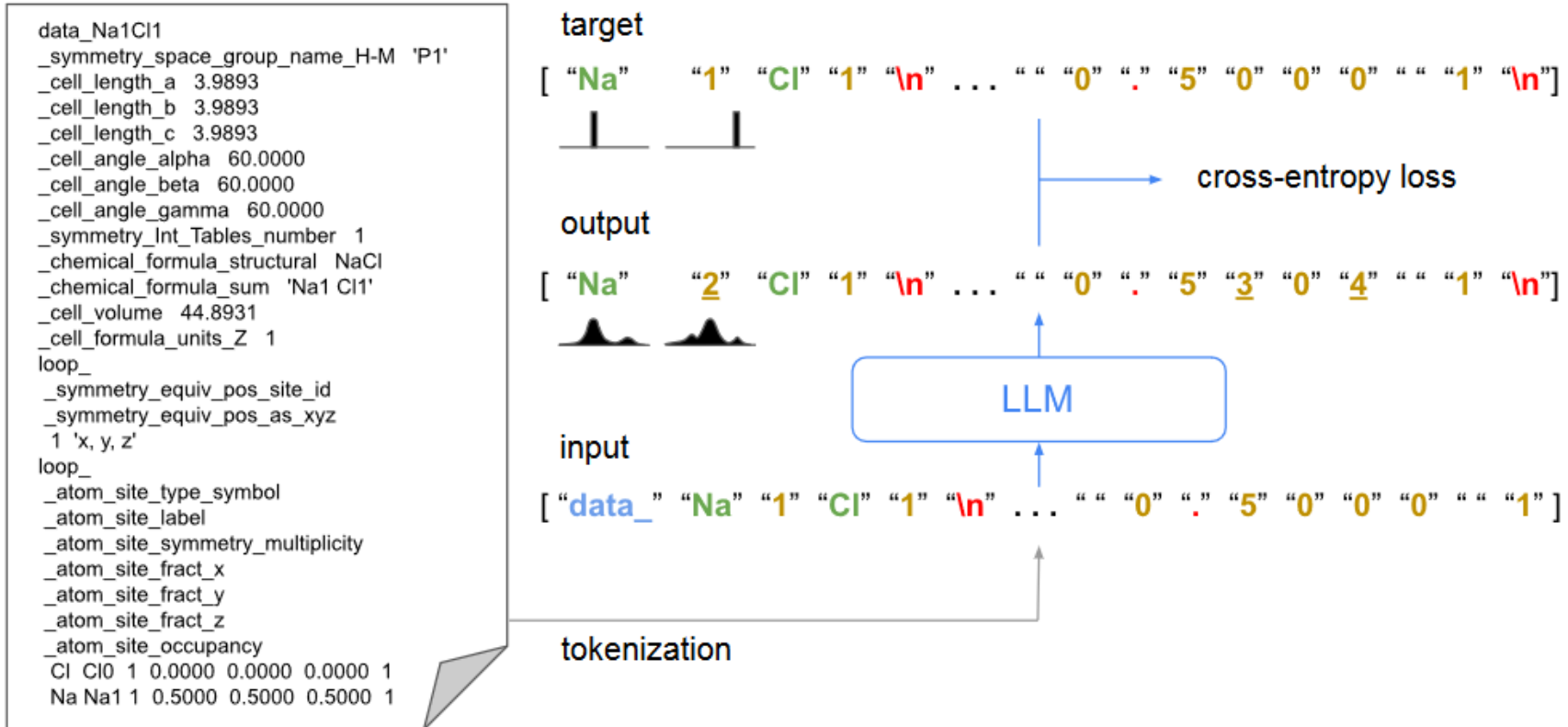
UPDATING THE EMBEDDINGS



$$E \longrightarrow \text{softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \longrightarrow E'$$

Each of these softmax matrix multiplications is a 'head'

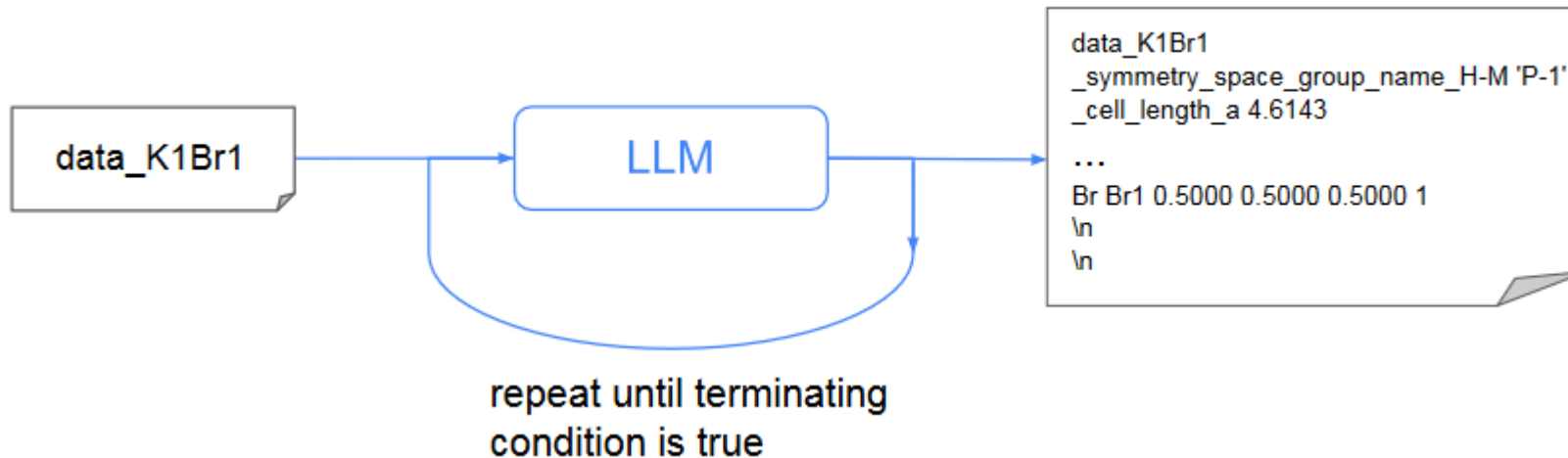
CRYSTALLM



A **decoder only transformer** trained on cif files for materials **structure generation**

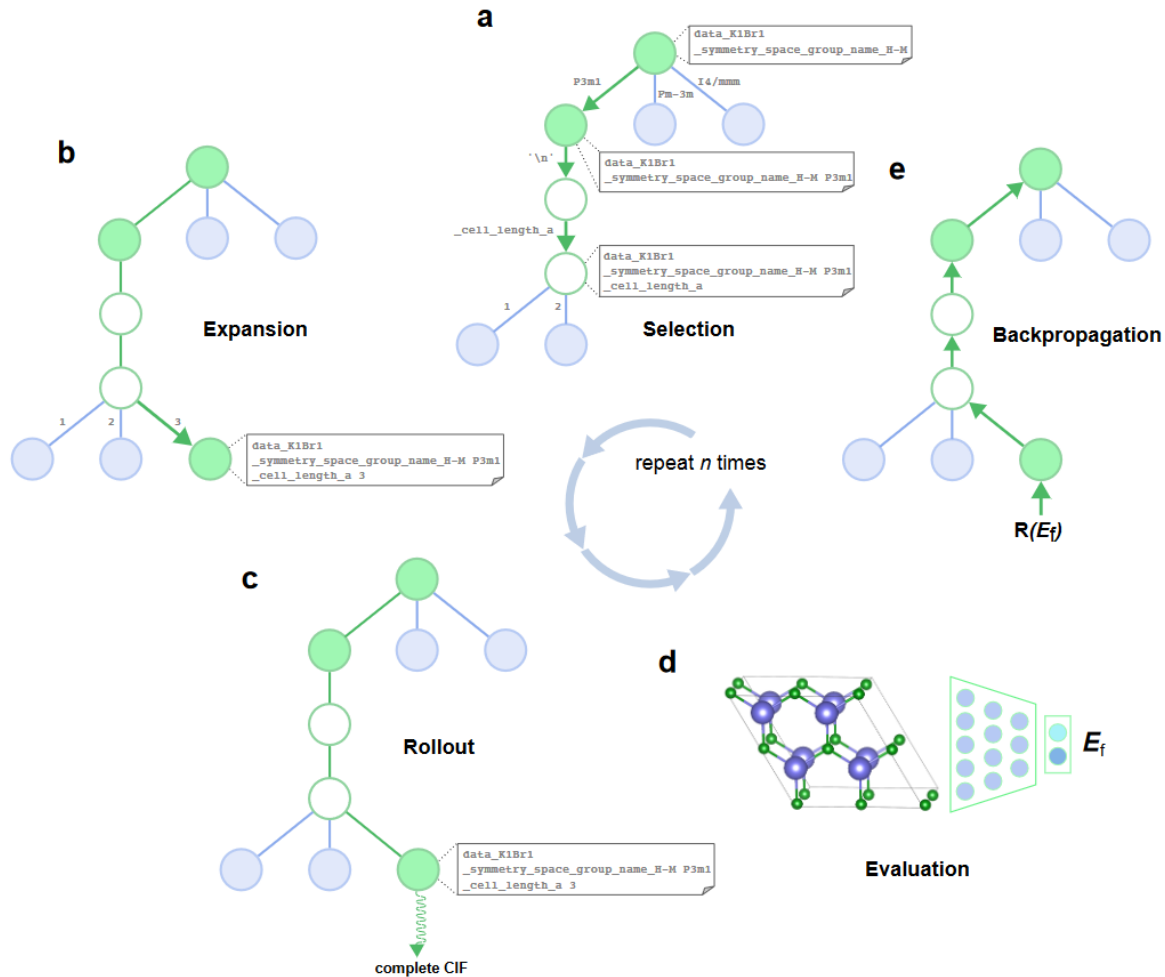
AUTOREGRESSIVE GENERATION

After training, CrystaLLM can be prompted with new text and can produce a predicted complete cif



Prompting is flexible so we can provide as little or as much information as we like

MONTE CARLO TREE SEARCH FOR CONSISTENCY



Autoregressive generation is stochastic and can lead to non-ideal structures

MCTS is more expensive but uses an energy estimator to drive to low energy solutions

CONCEPT CHECKLIST

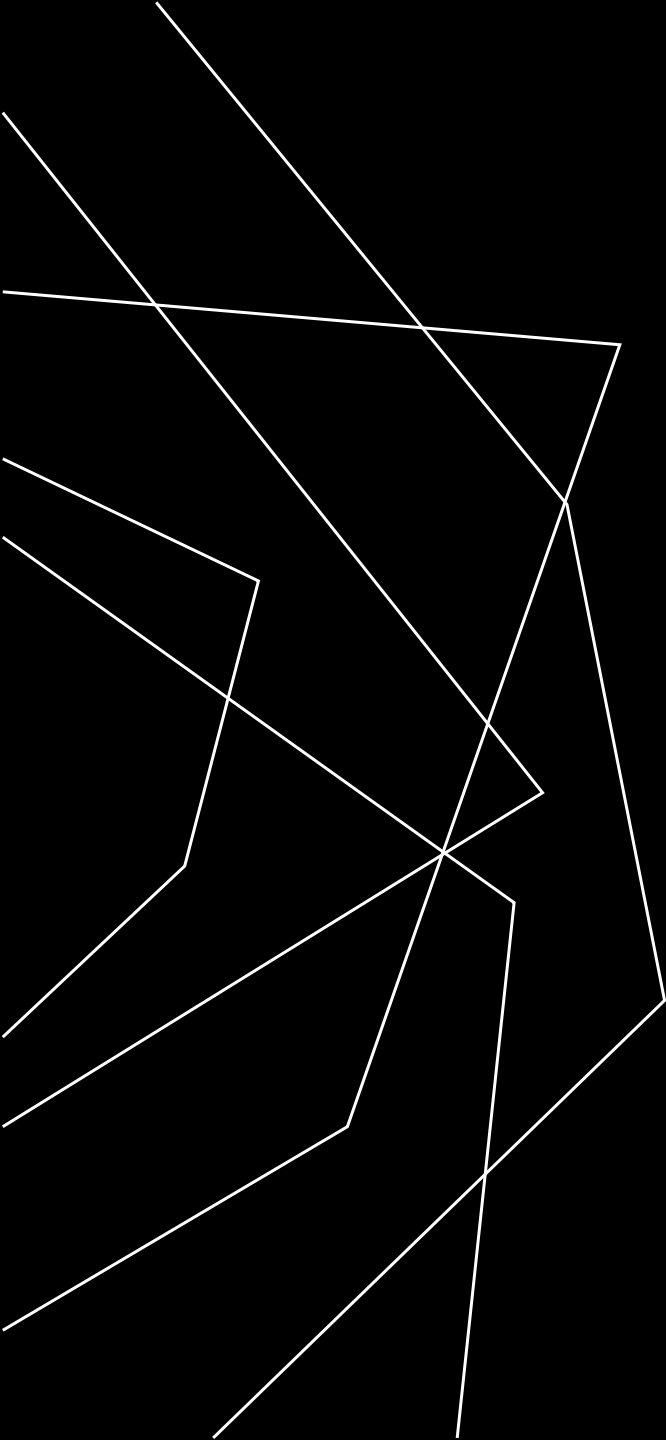
Sequential data benefits from **contextual awareness** and memory

Recurrent networks are an early answer

Memory was improved with LSTMs

Transformers use **attention** to map across sequences

Autoregression can be applied to generate **crystal structures**



THANK YOU

mdi-group.github.com