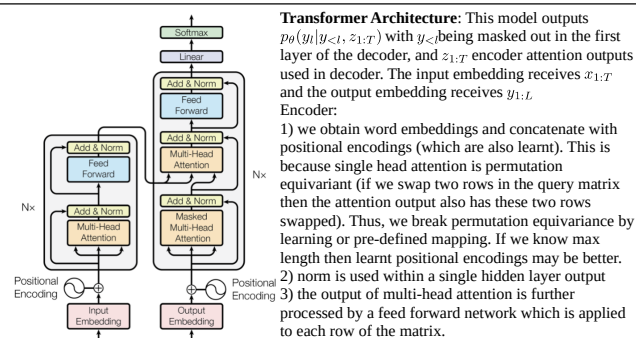
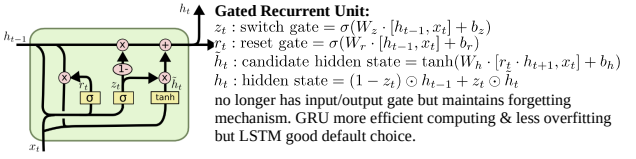
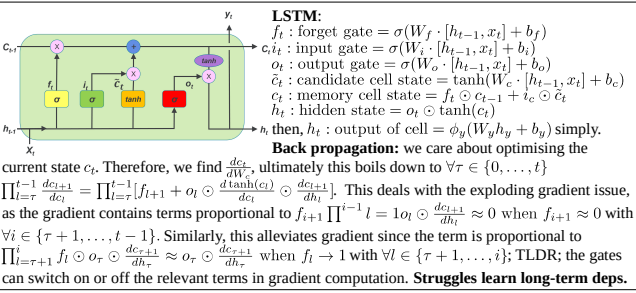
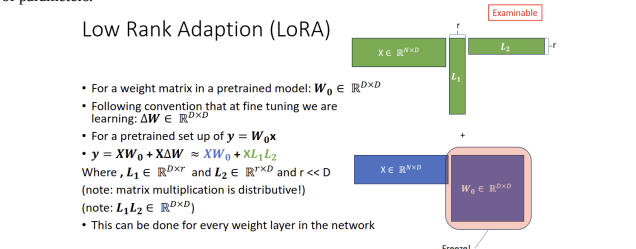


RNN: $h_t = \phi_h(W_h h_{t-1} + W_x x_t + b_h)$ and $y_t = \phi_t(W_y h_t + b_y)$ where ϕ is the activation.
Backward Pass: calculate $\mathcal{L}(\theta) = L_{total}(\theta) = \sum_{t=1}^T L(y_t)$ and $\frac{d}{d\theta} \mathcal{L}(\theta) = \sum_{t=1}^T \frac{d}{d\theta} L(y_t)$ is computed for each $\theta \in \{W_h, W_x, W_y, b_h, b_y\}$.
• $\frac{d\mathcal{L}(y_t)}{dW_x} = \frac{d\mathcal{L}(y_t)}{dy_t} \frac{dy_t}{dh_t} \frac{dh_t}{dW_x} = \frac{d\mathcal{L}(y_t)}{dh_t} \frac{dh_t}{dW_x}$
• $\frac{d\mathcal{L}(y_t)}{dW_h} = \frac{d\mathcal{L}(y_t)}{dy_t} \frac{dy_t}{dh_t} \frac{dh_t}{dW_h} = \frac{d\mathcal{L}(y_t)}{dh_t} \frac{dh_t}{dW_h}$. However, W_x, b_h contribute to h_t directly and indirectly. Next bullet-point describes why we also need $\frac{dh_t}{dW_x} = \frac{\partial h_t}{\partial W_x} + \frac{dh_{t-1}}{dW_x}$ and $\frac{dh_t}{db_h} = \frac{\partial h_t}{\partial b_h} + \frac{dh_{t-1}}{db_h}$. The entries in the Jacobian $\frac{dh_t}{dW_x}$ contains the total gradient of $h_t[i]$ w.r.t $W_h[m, n]$. h_t depends on h_{t-1} which depends W_h . Therefore: $\frac{dh_t}{dW_h} = \frac{\partial h_t}{\partial W_h} + \frac{dh_{t-1}}{dW_h} \frac{dh_{t-1}}{dW_h}$
• $\frac{d\mathcal{L}(y_t)}{dW_x} = \frac{d\mathcal{L}(y_t)}{dy_t} \frac{dy_t}{dh_t} \frac{dh_t}{dW_x} = \frac{d\mathcal{L}(y_t)}{dh_t} \frac{dh_t}{dW_x}$. We can truncate this to $\sum_{\tau=1}^{t-1} (\prod_{l=\tau}^{t-1} \frac{dh_{l+1}}{dh_l}) \frac{\partial h_{l+1}}{\partial W_x}$ which results in the **Back-propagation through time**. We may truncate this to $\sum_{\tau=\max(1, t-L)}^t$ \prod contains products of activation and weight matrix; gradient vanish/explode as $t \rightarrow \infty$!
Therefore the long-term dependencies become harder to learn.



Decoder:
1) after embedding and positional encoding, masked attention is used so during training only the previous words can be used.
2) The input representation is fed into the multi-head attention. The encoder output $z_{1:T}$ is used as the keys and values of the attention module. This allows the decoder to attend every word in the input $x_{1:T}$ for each of the predicted outputs $y_{1:L-1}$ so far.
The network is trained using maximum likelihood or cross entropy loss.

Fine Tuning: if you have a pre-trained model you want to translate it into another task.
If you have pre-trained weights ϕ , you need to update $\phi + \Delta\phi$. We have an auto-regressive model P_ϕ and context target tokens (that we're fine tuning on) $\max_\phi \sum_{\{x, y\} \sim Z} \sum_{t=1}^y \log(P_\phi(y_t | x, y_{<t}))$ this requires us to store info about every parameter ϕ . This is also bad if you have a lot of fine-tuned downstream tasks.
Low Rank Adaption (LoRA): what if $\Delta\phi$ also has low intrinsic rank we can approximate $\Delta\phi$ with Θ where $|\Theta| < |\phi|$. $\max_\Theta \sum_{\{x, y\} \sim Z} \sum_{t=1}^y \log P_{\phi+\Theta}(y_t | x, y_{<t})$ i.e. only update a subset of parameters.



Surprisingly, $r=1,2$ give good results. At inference time you combine $W_0 + L_1L_2$ and we can swap L_1L_2 s. However, you still need to store the parameters in order to do inference.
QLoRA: this is going from a data type of higher precision or a higher bit rate to a lower bit rate. $X^{Int8} = \text{round}(\frac{1}{2^{FP32}} X^{FP32}) = \text{round}(e^{FP32} \cdot X^{FP32})$ and for dequantising: $\text{dequant}(e^{FP32}, X^{Int8}) = \frac{X^{Int8}}{e^{FP32}} = X^{FP32}$ This is lossy. If you have outliers, the representation becomes sparse. Thus chunk the data and quantize each block separately with its own c^{*x}
QLoRA: BrainFloat: deep learning is more sensitive to underflow than overflow.
QLoRA: Quantising W_0 from bfloat16 to 4 bit: store weights with BrainFloat in 4 bits
Double Quantisation: quantise the quantisation constants.
Finally, for inference, dequantise with $y = X \text{dequant}(W_0) + XL_1L_2$

Sequence to Sequence (Seq2Seq) Auto-regressive Model: given two varying sentence length pairs, $p_\theta(y_{1:L} | x_{1:T}) = \prod_{t=1}^L p_\theta(y_t | y_{<t}, v)$, $v = \text{enc}(x_{1:T})$ fit an LSTM model 'sequence decoder' with $p_\theta(y_t | y_{<t}, v) = p_\theta(y_t | h_t^d, c_t^d)$, $h_t^d, c_t^d = \text{LSTM}_\theta^{dec}(y_{<t})$ where the first hidden states h_0^d, c_0^d are initialised by an encoder $v = \text{enc}(x_{1:T}) = N N_\theta(h_T^e, c_T^e)$, $h_T^e, c_T^e = \text{LSTM}_\theta^{enc}(x_{1:T})$. We initialise the system by encoding the entirety of the x input. This LSTM will return the first word (upon seeing the <eos> and use the hidden state v to initialise the decoder LSTM. Alternatively, the first word can be an <sos> token. **This can be extended to image captioning** where we extract an image encoding with CNNs and feed it into an LSTM.
Generative Models for Sequence VAE: We use an encoder architecture like an LSTM to predict the mean and standard variation of the q distribution. The loss function is defined similarly to before:
$$\phi^*, \theta^* = \arg \max \mathcal{L}(\phi, \theta), \quad \mathcal{L}(\phi, \theta) = \mathbb{E}_{p_{data}(x_{1:T})} [\underbrace{\mathbb{E}_{q_\phi(z|x_{1:T})} [\log p_\theta(z|x_{1:T})] - \text{KL}[q_\phi(z|x_{1:T}) || p_\theta(z)]}_{:= \mathcal{L}(x_{1:T}, \phi, \theta)}]$$

Encoder: $q_\phi(z|x_{1:T}) = N(z; \mu_\phi(x_{1:T}), \text{diag}(\sigma_\phi^2(x_{1:T})))$, where $\mu_\phi(x_{1:T}), \log \sigma_\phi(x_{1:T}) = \text{LSTM}_\phi(x_{1:T})$ and the Generator: $p_\theta(x_{1:T}|z) = \prod_{t=1}^T p_\theta(x_t | x_{<t}, z)$.

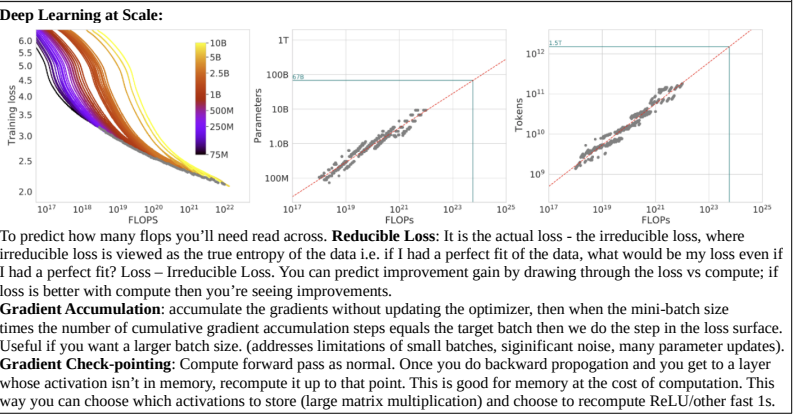
State-Space-Model: stochastic dynamic model for the latent space $p_\theta(z_{1:T}) = p_\theta(z_1) \prod_{t=2}^T p_\theta(z_t | z_{t-1})$. The emission depends on z_t only: $p_\theta(x_t | z_t)$. The joint distribution is thus $p_\theta(x_{1:T}, z_{1:T}) = \prod_{t=1}^T p_\theta(z_t | z_{t-1}) p_\theta(x_t | z_t)$
Hidden Markov Model (HMM) $z_t = A z_{t-1} + B \epsilon_t, \epsilon_t \sim N(0, I)$ with a Linear Gaussian emission model $x_t = C z_t + D \psi_t, \psi_t \sim \epsilon_t$. **State-space models + non-linear dynamics:** embed an LSTM which will predict the mean and variance for the next latent variable z at $t+1$: $p_\theta(y_t | z_t) = \mu_\theta^z(z_t) + \sigma_\theta^z(z_t) \psi_t, \psi_t \sim N(0, I)$ with the joint distribution $p_\theta(x_{1:T}, z_{1:T}) = \prod_{t=1}^T p_\theta(z_t | z_{t-1}) p_\theta(x_t | z_t)$. We train it with
$$E_{p_{data}(x_{1:T})} [\log p_\theta(x_{1:T})] \geq E_{p_{data}(x_{1:T})} [E_{q_\phi(z_{1:T}|x_{1:T})} [\log p_\theta(x_{1:T}|z_{1:T})] - \text{KL}[q_\phi(z_{1:T}|x_{1:T}) || p_\theta(z_{1:T})]]$$

For the encoder, we can base the q distribution based on the current input x with: $q_\phi(z_{1:T}|x_{1:T}) = \prod_{t=1}^T q_\phi(z_t | x_{\leq t})$ where $q_\phi(z_t | x_{\leq t}) = N(z_t; \mu_\phi^z(h_t^e), \text{diag}(\sigma_\phi^z(h_t^e)^2))$, $[h_t^e, c_t^e] = \text{LSTM}_\theta(x_t, h_{t-1}^e, c_{t-1}^e)$

Attention: before, we use v (the hidden state from the last LSTM block). How, use v_t such that with attention: $p_\theta(y_{1:L} | x_{1:T}) = \prod_{t=1}^L p_\theta(y_t | y_{<t}, v_t)$ with $v_t = \sum_{i=1}^T \alpha_{it} f_i$ which are the features each weighted by attention α_{it} and $h_t^d, c_t^d = \text{LSTM}_\theta^{dec}([y_{<t}, v_t])$, h_{t-1}^d, c_{t-1}^d where $f_t = T_\theta(h_t^e)$ (feature output of the encoder at time t) and similarity $e_t = a(h_{t-1}^e, f_t)$ and $\alpha_t = \text{softmax}(e_t)$, $e_t = \{e_{t1}, \dots, e_{tT}\}$
Single Head Attention: each key is paired with a value in the V matrix. The query V matrix contains query matrices. **Hard attention** is when $a = \text{onehot}(\arg \max x_i)$ and **Soft attention** is when $a = \text{softmax}(x)$ **Masked attention** is when we mask out some attention values (mask is 0 or 1 and sets the cell in the QK^\top matrix as $-\infty$ (useful for sequence prediction with a given ordering: in test time "future" is not available for the "current" to attend).
Complexity: time $O(MN d_q + MN d_v)$, space $O(MN + N d_v)$, # params (only key and value) so $K, V : O(M d_q + M d_v)$ and in self attention this is $O(N d_v)$ i.e. $Q=V=K$

Multi Head Attention: $\text{Multihead}(Q, K, V, a) = \text{concat}(\text{head}_1, \dots, \text{head}_N) W^O$, where $\text{head}_i = \text{Attention}(QW_i^K, KW_i^K, VW_i^V, a)$ and Attention is defined as before. This helps define different hidden similarity measures. **Complexity:** $W_i^K \in \mathbb{R}^{d_q \times d_q}, W_i^V \in \mathbb{R}^{d_v \times d_v}, W_i^O \in \mathbb{R}^{h d_v \times d_{out}}$ then time complexity $O(h(MN \tilde{d}_q + MN \tilde{d}_v) + h(\tilde{d}_q d_q (M + N) + \tilde{d}_v d_v M) + N \tilde{h} \tilde{d}_v d_{out})$
space complexity $O(h(MN + N \tilde{d}_v) + h(\tilde{d}_q (M + N) + \tilde{d}_v M) + N \tilde{d}_{out})$ where to keep complexity similar we set $\tilde{d}_q = \lfloor \frac{d_q}{h} \rfloor, \tilde{d}_v = \lfloor \frac{d_v}{h} \rfloor$ to keep the costs close to performing single-head attention in the original space.

Floating Point is a way to represent numbers in a computer where the decimal can float. We have a sign, exponent and fraction. $\text{value} = (-1)^{\text{sign}} \times 2^{E - 2^{\lfloor \text{exponent} \rfloor - 1} + 1} \times (1 + \sum_{i=1}^{\text{remaining}} \text{value} \times 2^{-i})$.
Calculate Memory footprint: Inference: bits per parameter \times number parameters. **Training:** Model params (same), Gradients: bits per gradient \times number parameters (gradient for each) **Optimizer:** need momentum and rms prop value: $2 \times$ same again. **Activations** also.
Floating Point Operations (FLOP): for $w, x \in \mathbb{R}^n, w + x : n \text{ FLOPs}, w^\top x : 2n \text{ FLOPs}$, for $W \in \mathbb{R}^{h \times \tilde{p}}, X \in \mathbb{R}^{\tilde{p} \times n} W X : 2n \tilde{p} \text{ FLOPs}$



To predict how many flops you'll need read across. **Reducible Loss:** It is the actual loss - the irreducible loss, where irreducible loss is viewed as the true entropy of the data i.e. if I had a perfect fit of the data, what would be my loss even if I had a perfect fit? Loss - Irreducible Loss. You can predict improvement gain by drawing through the loss vs compute; if loss is better with compute then you're seeing improvements.
Gradient Accumulation: accumulate the gradients without updating the optimizer, then when the mini-batch size times the number of cumulative gradient accumulation steps equals the target batch then we do the step in the loss surface. Useful if you want a larger batch size. (addresses limitations of small batches, significant noise, many parameter updates).
Gradient Check-pointing: Compute forward pass as normal. Once you do backward propagation and you get to a layer whose activation isn't in memory, recompute it up to that point. This is good for memory at the cost of computation. This way you can choose which activations to store (large matrix multiplication) and choose to recompute ReLU/other fast 1s.