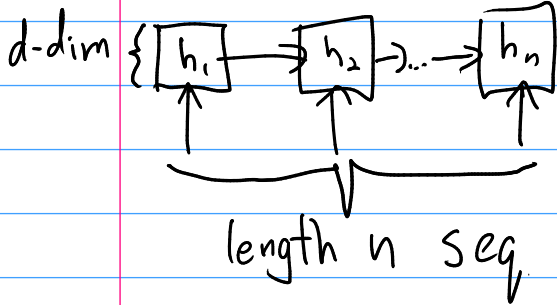


Attention

Before: RNNs for sequence model

$$h_2 = \sigma(W h_1 + Z x_2)$$

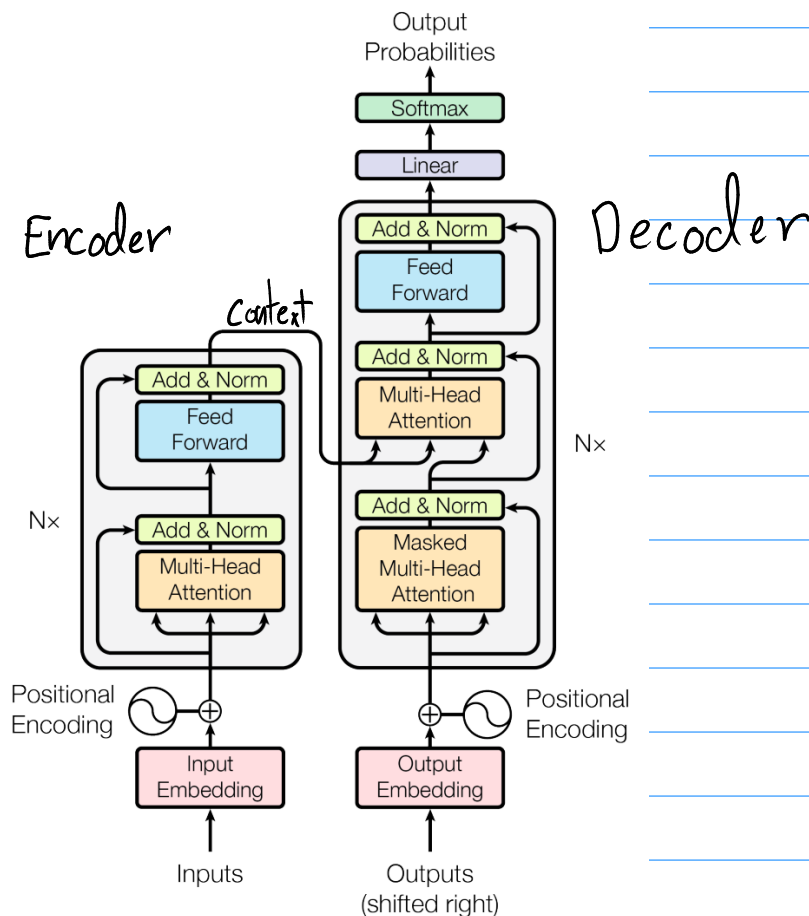


Computation: $O(d^2)$ per update
Total: $O(nd^2)$ computation

Not easily parallelizable
Longest path: $O(n)$ length
Some challenges in optimization

Attention mechanism

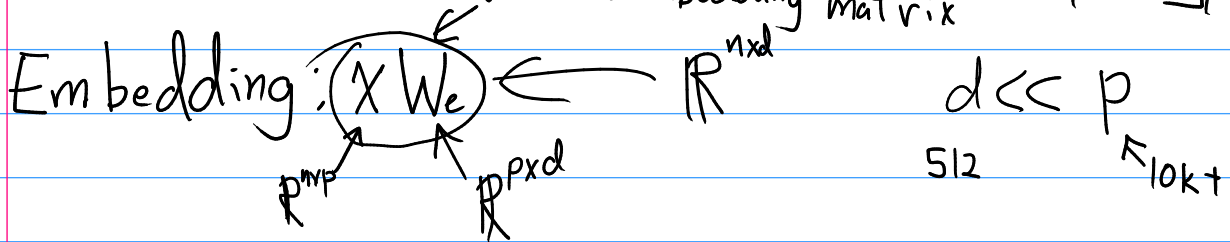
Transformer Architecture (Vaswani et al '17)



Input/output Encoding

$$X = (X_1, \dots, X_n)^T \in \mathbb{R}^{n \times p}, \quad X_i \in \mathbb{R}^p, \text{ one hot}$$

"The quick brown fox" \Rightarrow $\left(\begin{bmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ \dots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ \dots \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ \dots \\ 1 \end{bmatrix} \right)^T$



Positional encoding

The dog is behind the cat.
The cat is behind the dog.

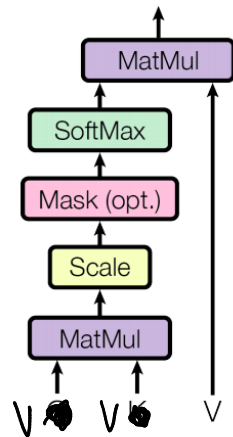
Construct $V \in \mathbb{R}^{n \times d}$:

a) Use fixed $V_{t, 2i} = \sin(t/10000^{2i/d})$ $i=0$ to $\frac{d}{2}-1$
 $V_{t, 2i+1} = \cos(t/10000^{2i/d})$

b) Learn a V

Add to embedding: $X W_e + V \in \mathbb{R}^{n \times d}$

Self

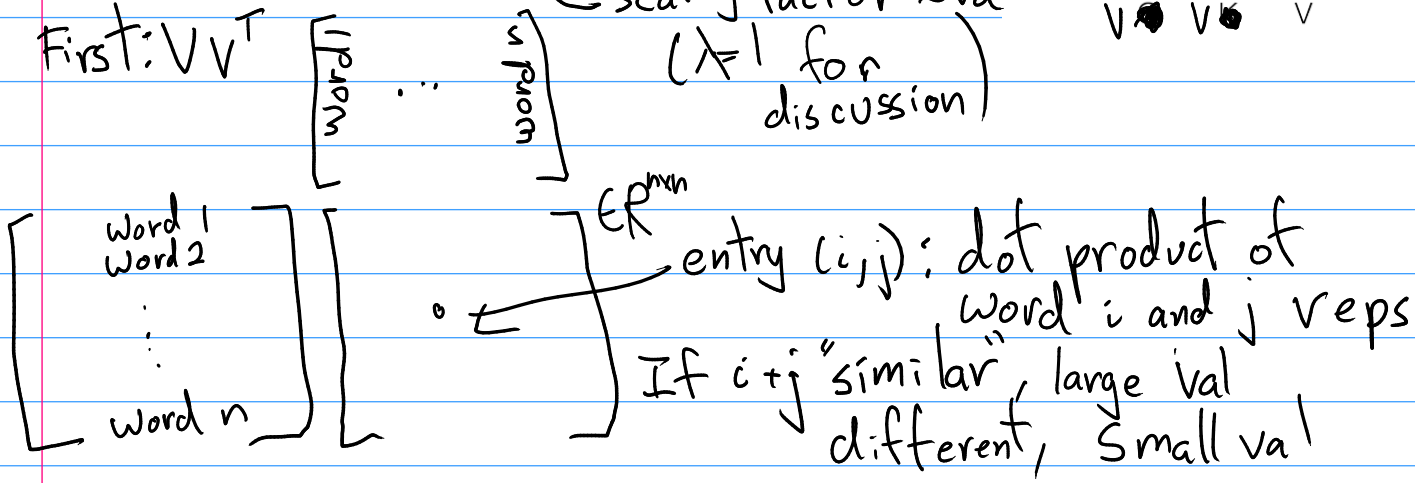


Self-Attention

Given $V \in \mathbb{R}^{n \times d}$, compute

$$V \leftarrow \text{softmax}(V V^T) V$$

First: $V V^T$ ↖ scaling factor $\approx \sqrt{d}$
($\neq 1$ for discussion)

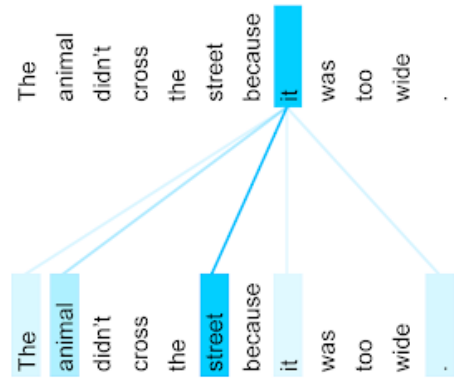
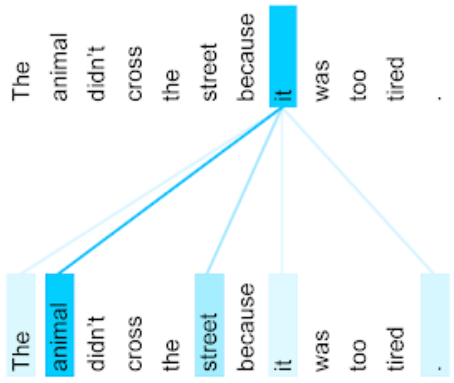


Softmax: make each row a distribution

$$\underbrace{\text{softmax}(V V^T)}_{\mathbb{R}^{n \times n}} \times \underbrace{V}_{\mathbb{R}^{n \times d}} = \mathbb{R}^{n \times d}$$

Replace each row of V by weight sum of rows

$$\begin{bmatrix} \frac{2}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{3}{4} \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & \frac{3}{4} & \frac{1}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{3}{4} & \frac{1}{2} & \frac{3}{4} \end{bmatrix}$$



Computation: $O(n^2d)$

$$VV^T \in O(n^2d)$$
$$\text{Softmax}(VV^T) V \in O(n^3d)$$

vs RNNs $O(nd^2)$

Shorter max path $O(1)$

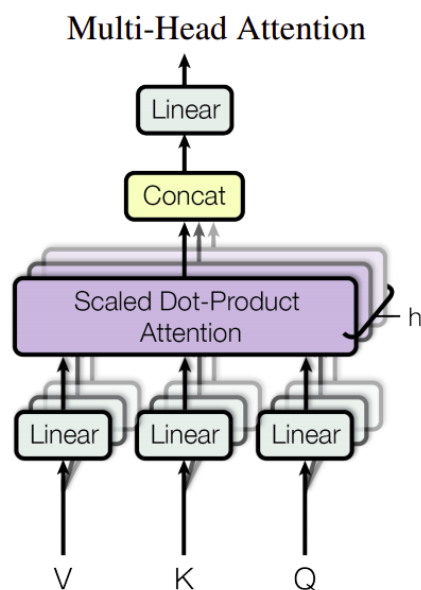
Multi-headed attention

For $i=1$ to h : $\mathbb{R}^{d \times \frac{d}{h}}$

$$V_i = V W_i^V$$

$$V_i = \text{softmax}(V_i V_i^T \Lambda) V_i$$

$$V \leftarrow \text{Concat}(V_1, \dots, V_h) W \quad \mathbb{R}^{pd \times d}$$



Other details

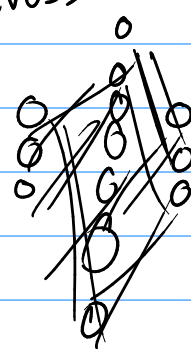
- Residuals, layer norm

- FF network

$$\text{FFN}(v) = W_2 \text{ReLU}(W_1 v)$$

Apply to each row of matrix

Shared across position



Decoder

Masked self attention

Given $Q \in \mathbb{R}^{n \times d}$, compute

$$Q \leftarrow \text{softmax}(\text{mask}(Q Q^T / \lambda)) Q$$

$$\text{mask}(M)_{ij} = \begin{cases} -\infty & \text{if } i < j \\ M_{ij} & \text{else} \end{cases}$$

$$M = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \quad \text{mask}(M) = \begin{bmatrix} 1 & -\infty & -\infty \\ 4 & 5 & -\infty \\ 7 & 8 & 9 \end{bmatrix}$$

Attention w/ context

Given context V , matrix Q , compute

$$Q \leftarrow \text{softmax}(Q^T V / \lambda) V$$

Final Softmax

Convert vectors from \mathbb{R}^d back to \mathbb{R}^p

gives pred seq \hat{y}

Optimize:

$$\min_{\text{(all params in model)}} \sum_{\text{all example}} \sum_{j=1}^l -y_j \log \hat{y}_j$$

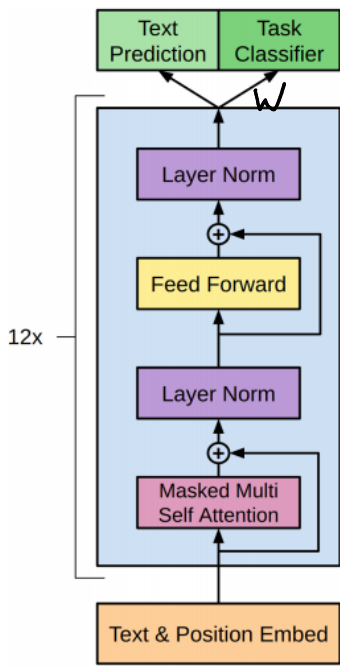
\leftarrow logit corresponding to true output

\uparrow true j th output

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

GPT-1

Two stages of training.



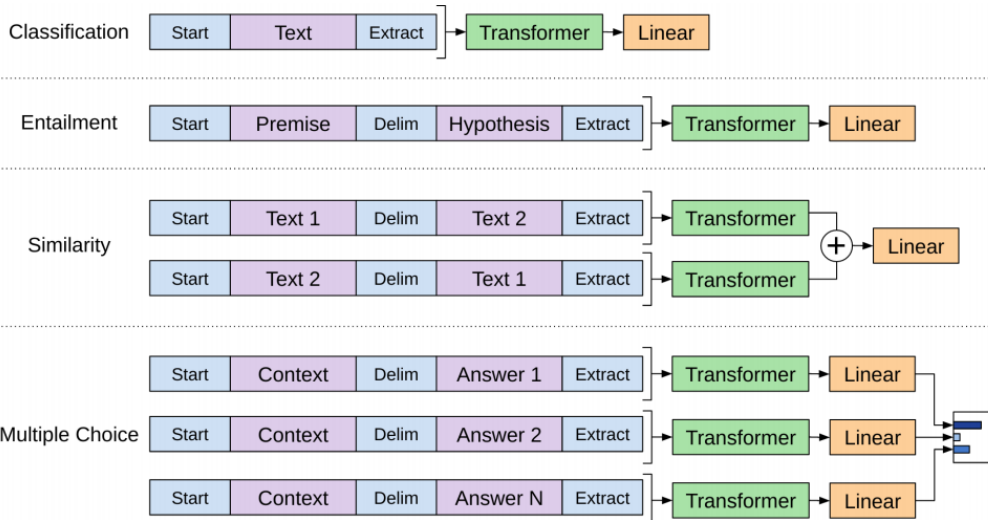
1. Unsupervised Pre-training
- lots of text used

- Objective: next word prediction

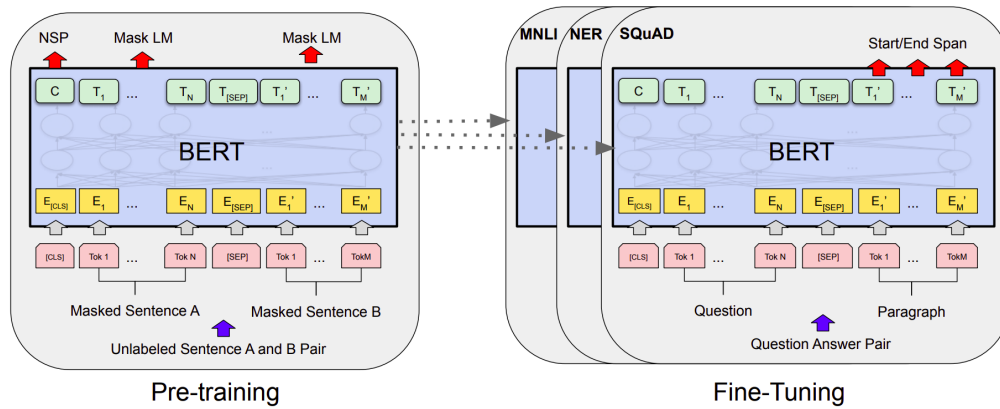
$$\min_{\theta} \sum_{\text{all examples}} \sum_{i=1}^n -\log P_{\theta}(x_i | x_1, \dots, x_{i-1})$$

2. Supervised Fine-tuning

$$\min_{\theta, W} \sum_{\text{all examples}} (-\log P_{\theta, W}(y|x)) + \lambda \sum_{i=1}^n -\log P_{\theta}(x_i | x_1, \dots, x_{i-1})$$



BERT



How to train:

a) Input: I took my [mask] for a walk.

⇒ I took my dog for a walk.

b) Feed in sentences A and B

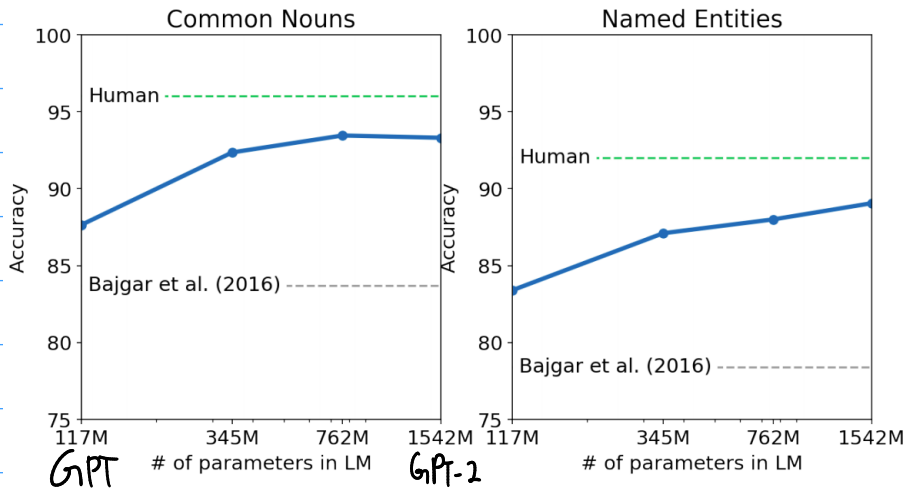
- Either B follows A, or is random

- Learn to predict which is the case

Sum two losses for pre training

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

GPT-2 10x bigger than GPT



GPT-3

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



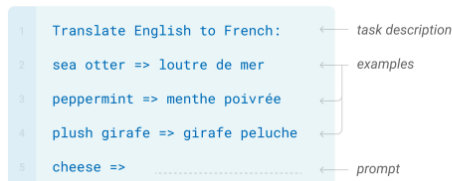
One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



Few-shot

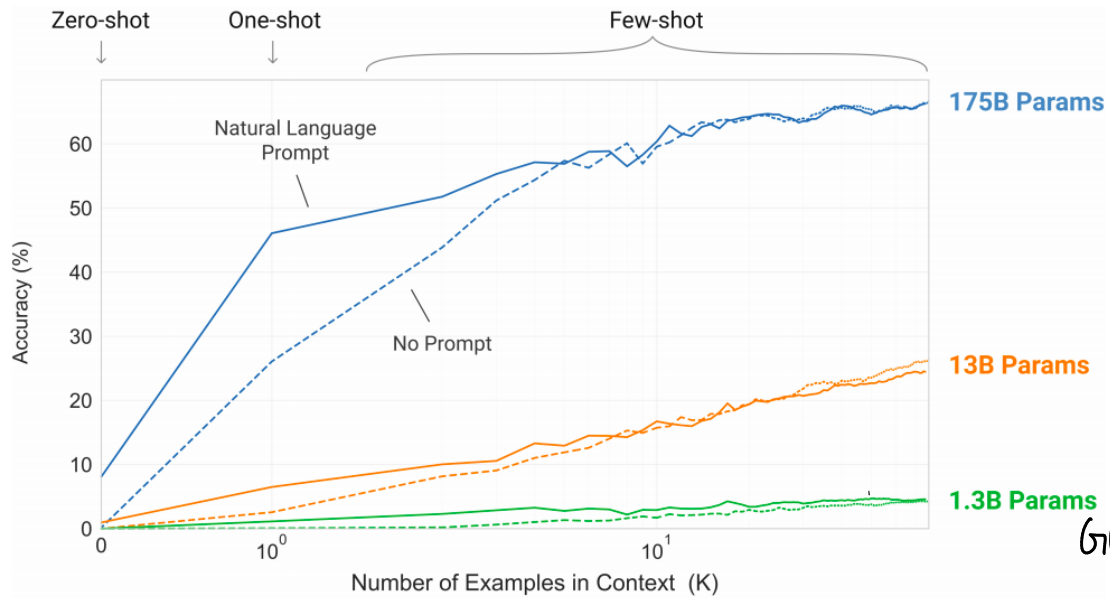
In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.





GPT-3

GPT-2

<p>Context →</p> <p>The City BY C. P. CAVAFY TRANSLATED BY EDMUND KEELEY [Poem text omitted]</p> <p>SOME TREES John Ashbery [Poem text omitted]</p> <p>Shadows on the Way Wallace Stevens</p>	<p>Generated Poem 1</p> <p>I must have shadows on the way If I am to walk I must have Each step taken slowly and alone To have it ready made</p> <p>And I must think in lines of grey To have dim thoughts to be my guide Must look on blue and green And never let my eye forget That color is my friend And purple must surround me too</p> <p>The yellow of the sun is no more Intrusive than the bluish snow That falls on all of us. I must have Grey thoughts and blue thoughts walk with me If I am to go away at all.</p> <p>Generated Poem 2</p> <p>There is a way in the world. That way a man may live in a small town His lifetime and see no man come, Never the stir of men about a market, Never the bobbing of covered carts, Never the creak passed on the stair Closing some other life; an old man Can prosper in that town, can prosper When his store is seen to close, his wife Fades in his house, and the world shouts in vain That no one comes to hear him Where his small bell tolls. In that small town Everything is known. The ways of the young men are known, And the ways of the young women As they walk to the Spring; The rumors are all known, The quarrels, the hurt, the envy Of the men and the women. Till a day comes That way out of the world That leads to a wood. Young, the young man Walks, a lover now, Knows the streets as they are and is of them.</p> <p>Comes to the wood.</p> <p>There, in the wood, among trees, He sees shadows on the way, hears voices, hears the wind and the rustling of leaves; Through an open glade He sees a shape and the shape hears: It waits as he waits, Waits as the shadows wait, As the voices wait; Shadows on the way, voices in the wind.</p>	<p>Generated Poem 3</p> <p>The sun was all we had. Now, in the shade All is changed. The mind must dwell on those White fields, that to its eyes were always old; Those ancient gleams, convoluting The way they lay among the huge roots, The great dark tomes of reverie, The plumed labyrinth of the sea.</p> <p>When our own breath gone faint and fade Up the long stairway of the dead We shall remember well The blindness of the light, the whiteness Of that white land. The footsteps, and The grasses in that light, and in the shade The menacing life.</p> <p>We played, too, once, in the warmest rooms. To one content, there is one content For all mankind. The forest glades Are all the more for being in fable, And wood and town in legend shadowed- Once more to return. The one content Has in it every happiness, more brave For being purposeless; and so Those island gardens, deep in the mind, Make of remoteness familiarity And of the unspoken name a friend.</p> <p>Generated Poem 4</p> <p>Nobody will come to this place. It is a road that leads nowhere. The solitude is deep. The mountains are high. But they are desolate, and they turn the traveler's face Towards the North. All the sounds of the world are far away. When the wind rises above the trees, The boughs bow to the ground. Even the birds that inhabit the tangle of weeds That is the roadside cover, are silent. One listens, But hears no roar of the forest. One is alone. One will be taken. One will be taken. There is no utterance, there is no conversation, But one is uneasy all the same.... There is a thin blue mist, A darkness rising like smoke, And within that darkness A possession of the heart. One will be taken... It was here, and it will be here again- Here, under this sky empty and full of light.</p>
---	---	---