

softmax	input:	(1, 16, 1370, 1370)
	output:	(1, 16, 1370, 1370)

Dropout	input:	(1, 16, 1370, 1370)
	output:	(1, 16, 1370, 1370)

matmul	input:	(1, 16, 1370, 1370), (1, 16, 1370, 64)
	output:	(1, 16, 1370, 64)

transpose	input:	(1, 16, 1370, 64)
	output:	(1, 1370, 16, 64)

reshape	input:	(1, 1370, 16, 64)
	output:	(1, 1370, 1024)

Linear	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

Dropout	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

LayerScale	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

add	input:	2 x (1, 1370, 1024)
	output:	(1, 1370, 1024)

LayerNorm	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

Linear	input:	(1, 1370, 1024)
	output:	(1, 1370, 4096)

gelu	input:	(1, 1370, 4096)
	output:	(1, 1370, 4096)

Dropout	input:	(1, 1370, 4096)
	output:	(1, 1370, 4096)

Linear	input:	(1, 1370, 4096)
	output:	(1, 1370, 1024)

Dropout	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

LayerScale	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

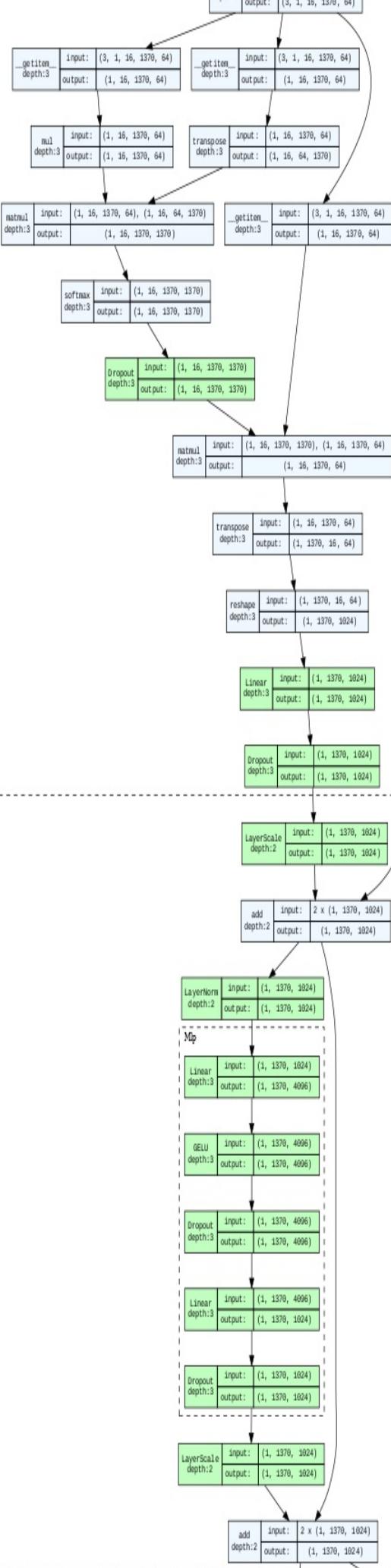
add	input:	2 x (1, 1370, 1024)
	output:	(1, 1370, 1024)

NestedTensorBlock	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

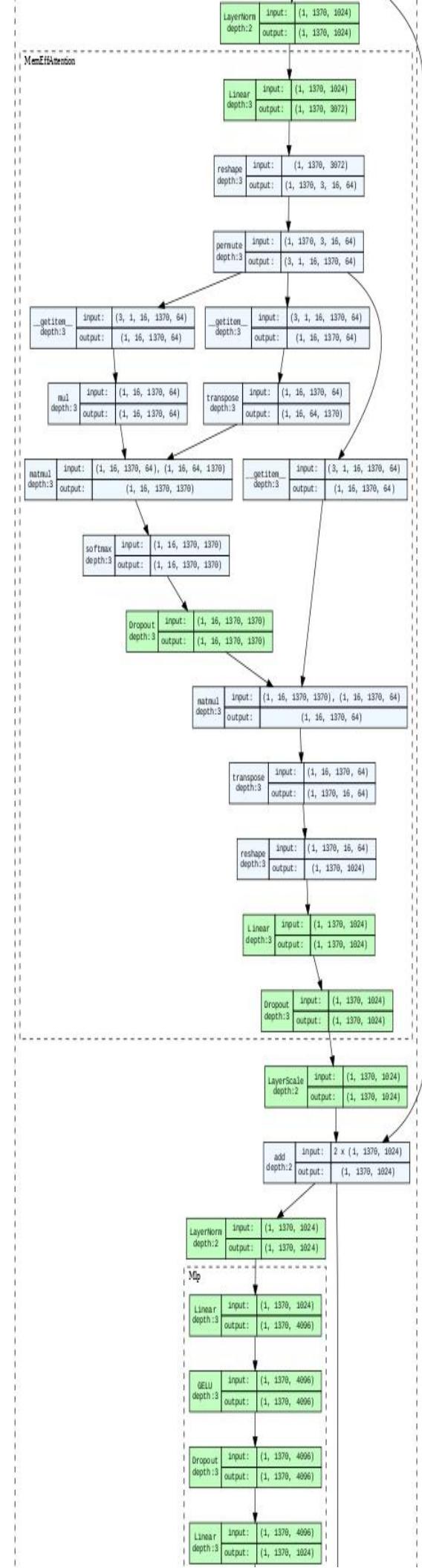
Linear	input:	(1, 1370, 1024)
	output:	(1, 1370, 3072)

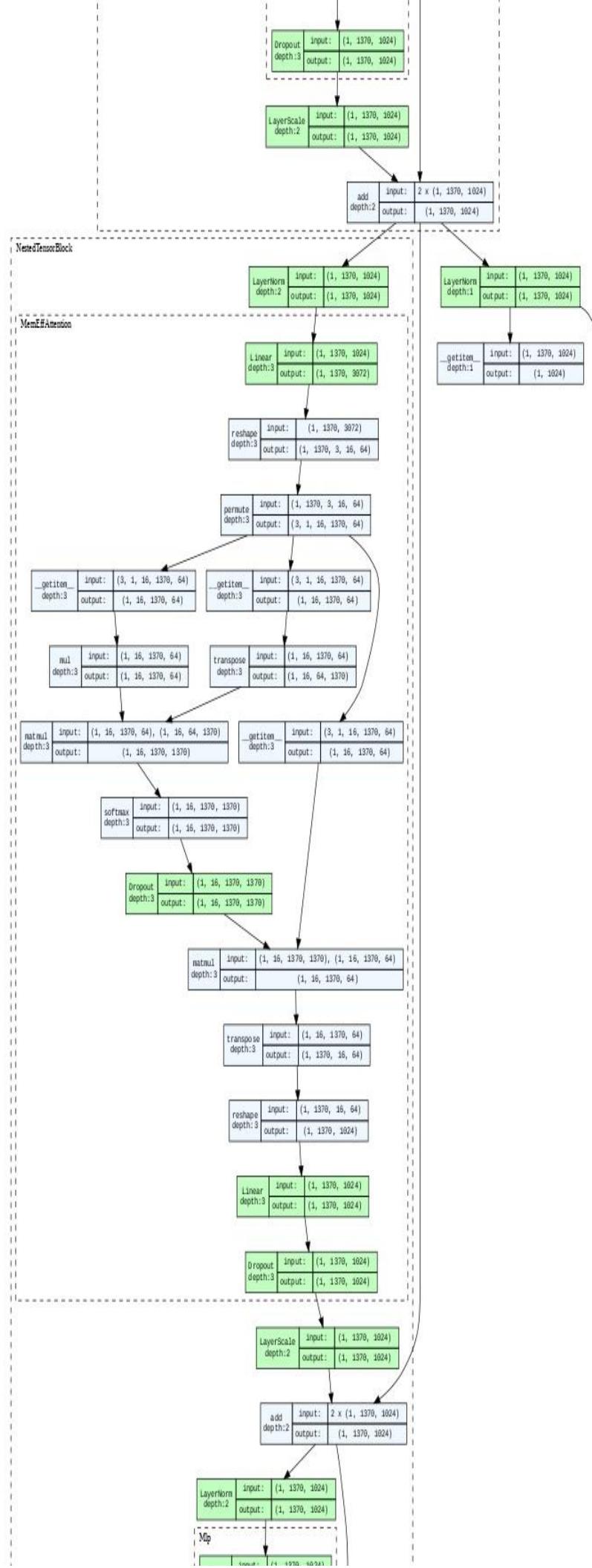
reshape	input:	(1, 1370, 3072)
	output:	(1, 1370, 3, 16, 64)

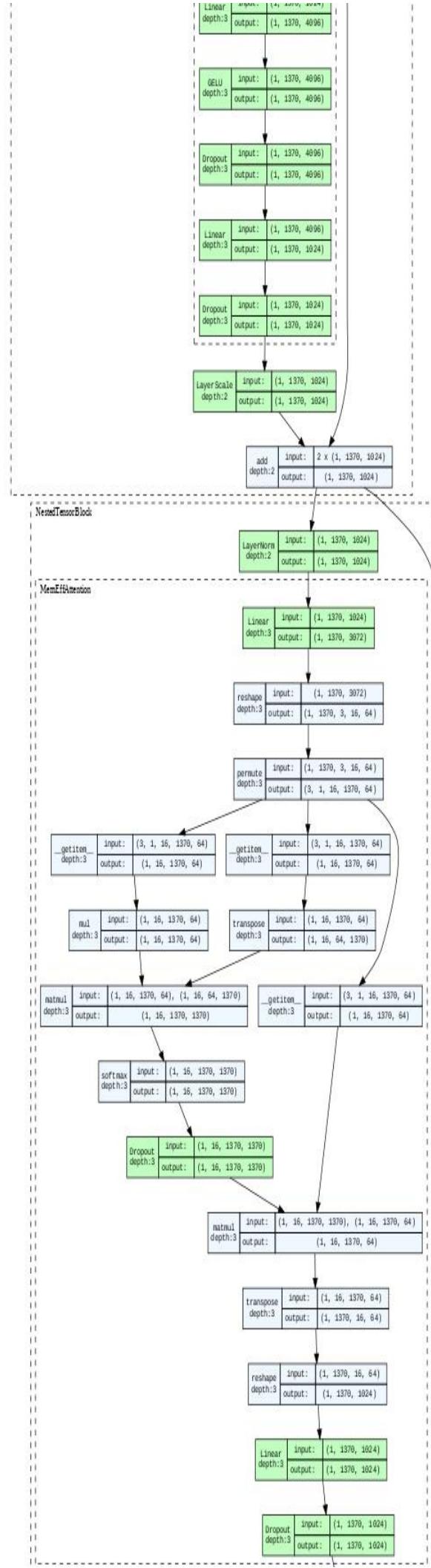
permute	input:	(1, 1370, 3, 16, 64)
	output:	(1, 1370, 3, 16, 64)

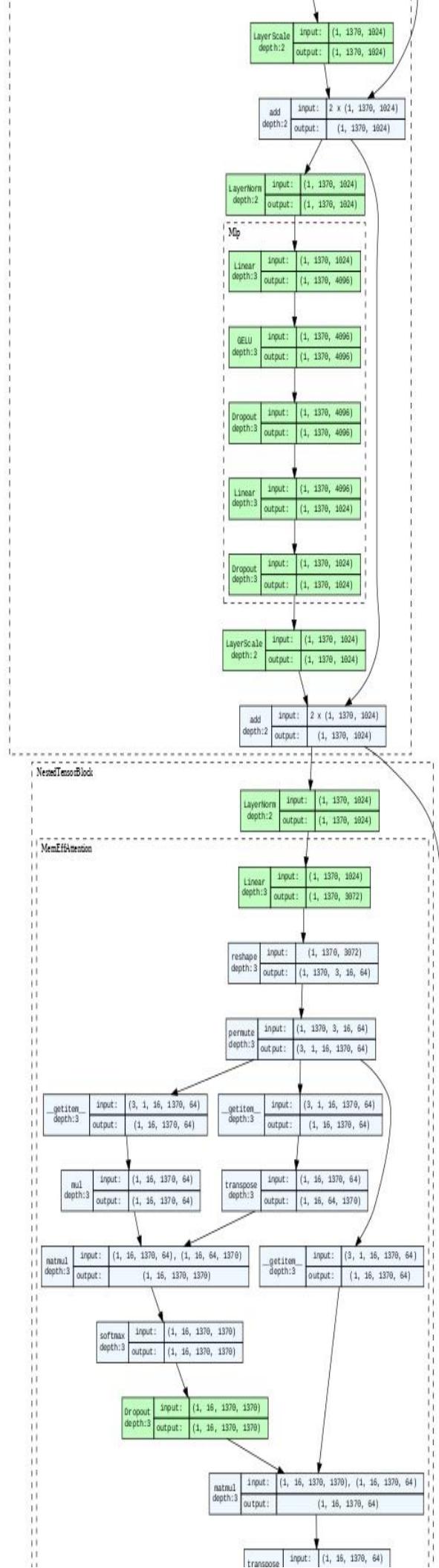


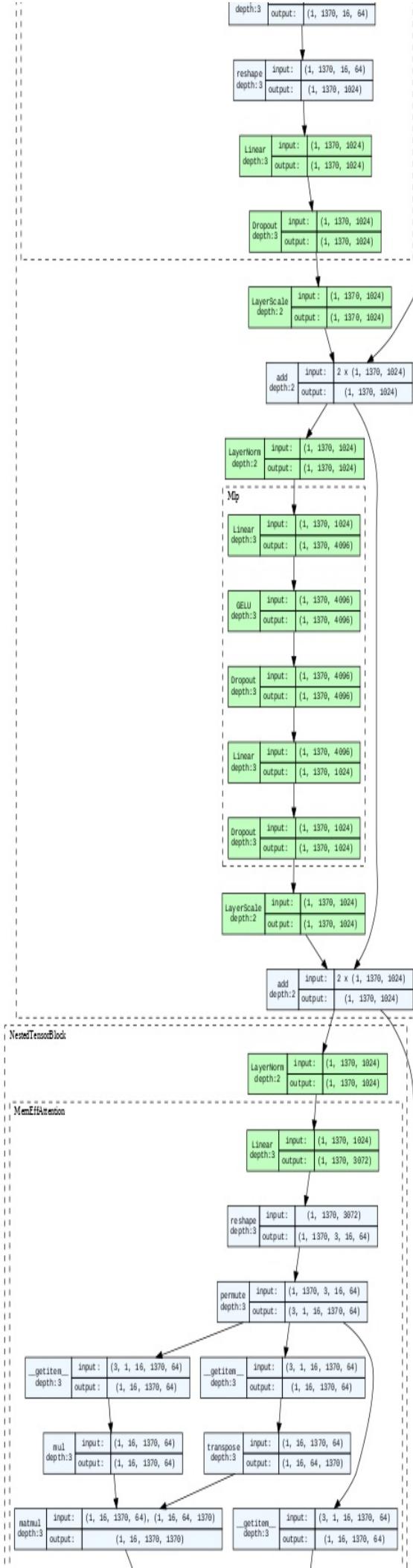
NestedTensorBlock











softmax	input:	(1, 16, 1370, 1370)
	output:	(1, 16, 1370, 1370)

Dropout	input:	(1, 16, 1370, 1370)
	output:	(1, 16, 1370, 1370)

matmul	input:	(1, 16, 1370), (1, 16, 1370, 64)
	output:	(1, 16, 1370, 64)

transpose	input:	(1, 16, 1370, 64)
	output:	(1, 1370, 16, 64)

reshape	input:	(1, 1370, 16, 64)
	output:	(1, 1370, 1024)

Linear	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

Dropout	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

LayerScale	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

add	input:	2 x (1, 1370, 1024)
	output:	(1, 1370, 1024)

LayerNorm	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

Mq	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

Linear	input:	(1, 1370, 1024)
	output:	(1, 1370, 4096)

GELU	input:	(1, 1370, 4096)
	output:	(1, 1370, 4096)

Dropout	input:	(1, 1370, 4096)
	output:	(1, 1370, 4096)

Linear	input:	(1, 1370, 4096)
	output:	(1, 1370, 1024)

Dropout	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

LayerScale	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

add	input:	2 x (1, 1370, 1024)
	output:	(1, 1370, 1024)

NestedTensorBlock

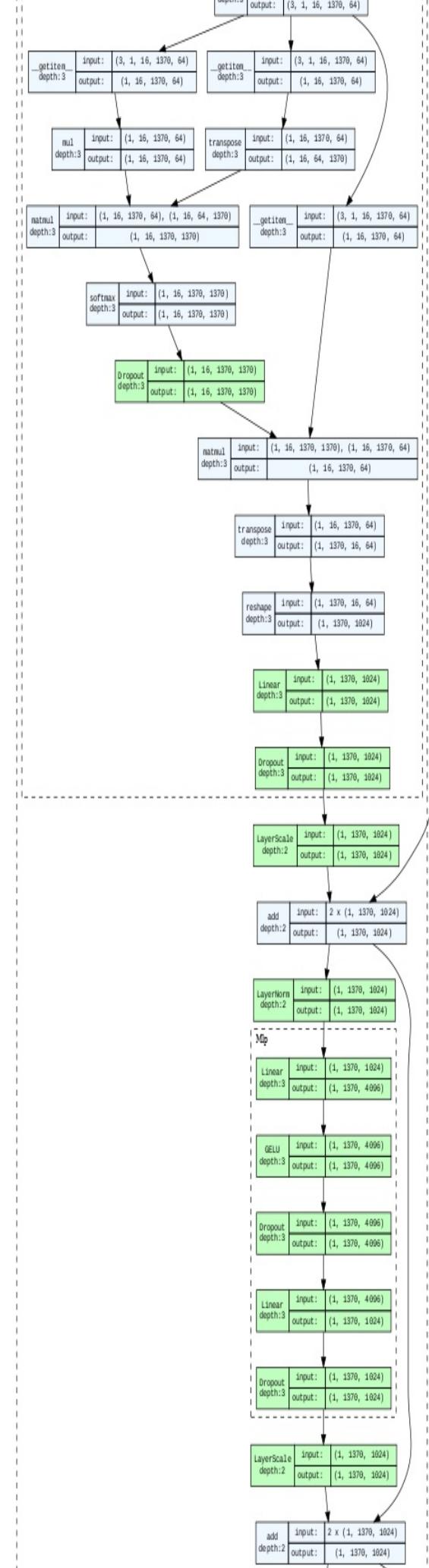
LayerNorm	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

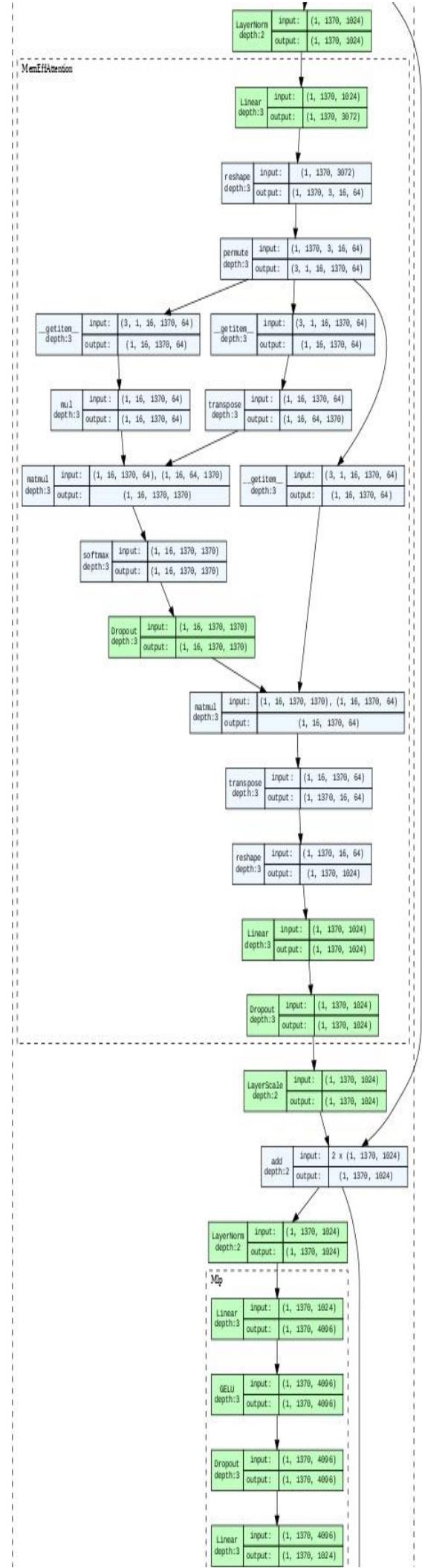
MemEffAttention

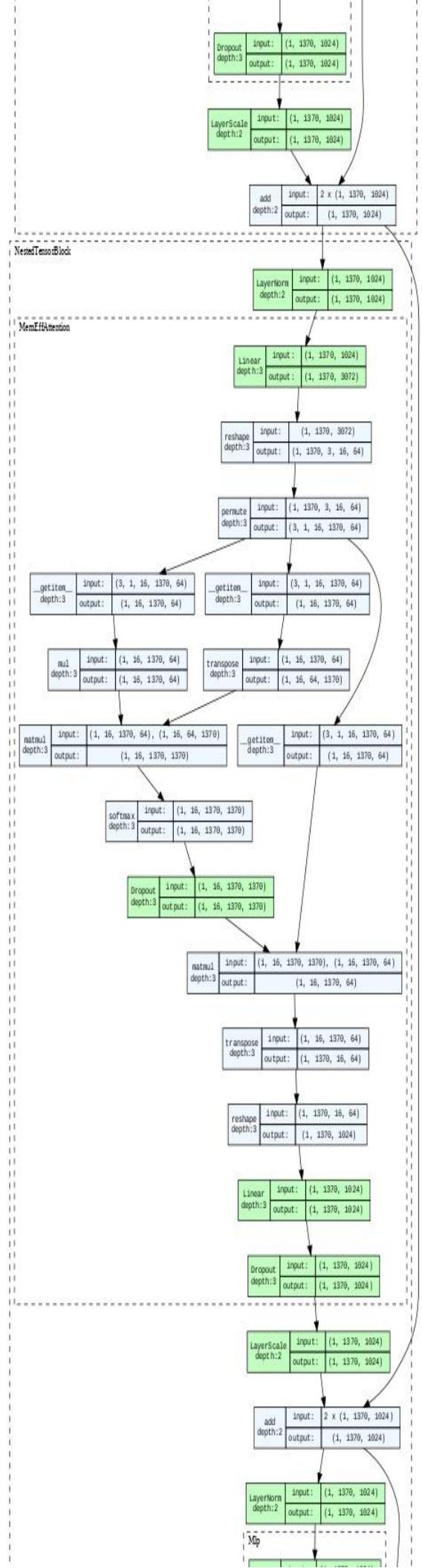
Linear	input:	(1, 1370, 1024)
	output:	(1, 1370, 3072)

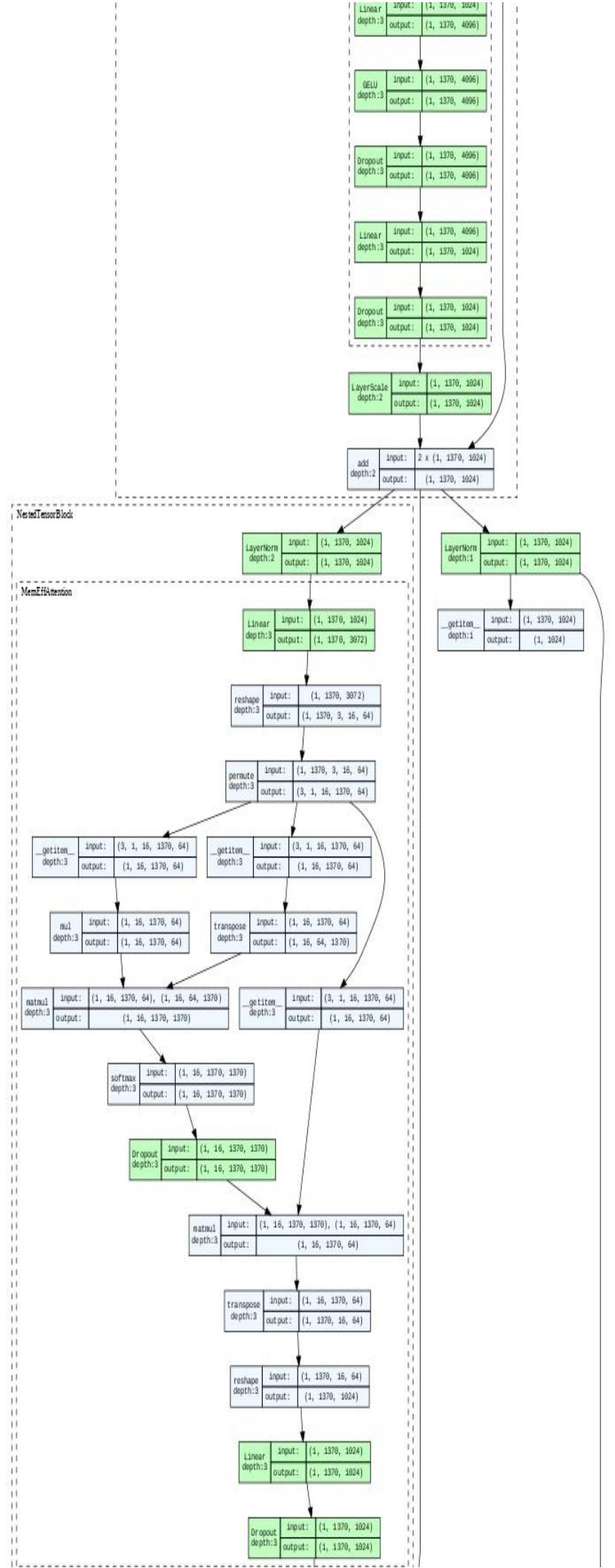
reshape	input:	(1, 1370, 3072)
	output:	(1, 1370, 3, 16, 64)

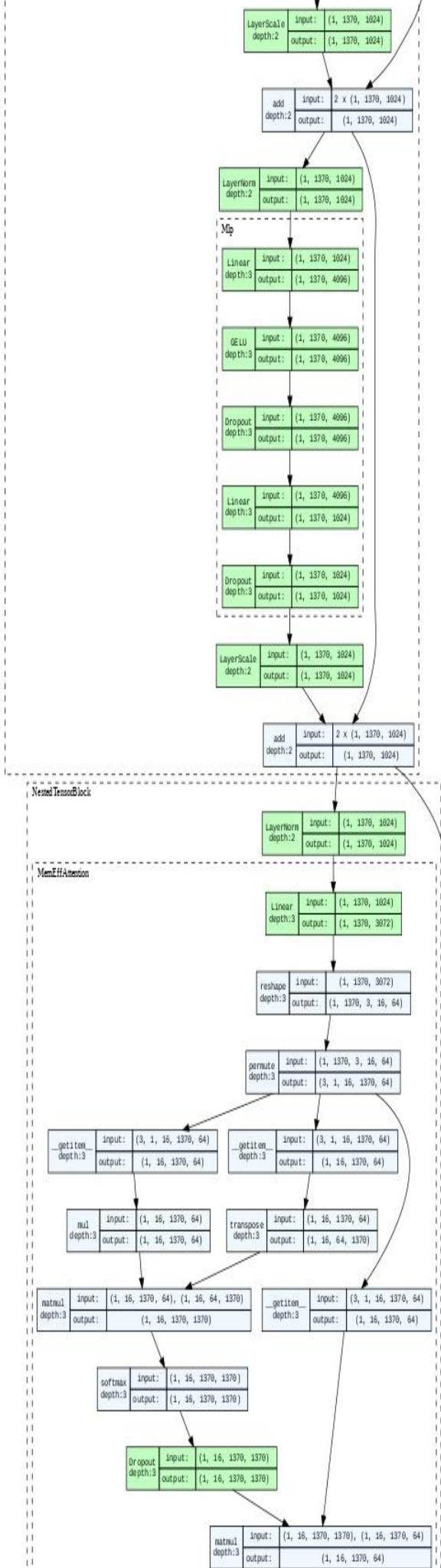
permute	input:	(1, 1370, 3, 16, 64)
	output:	(1, 1370, 3, 64, 16)

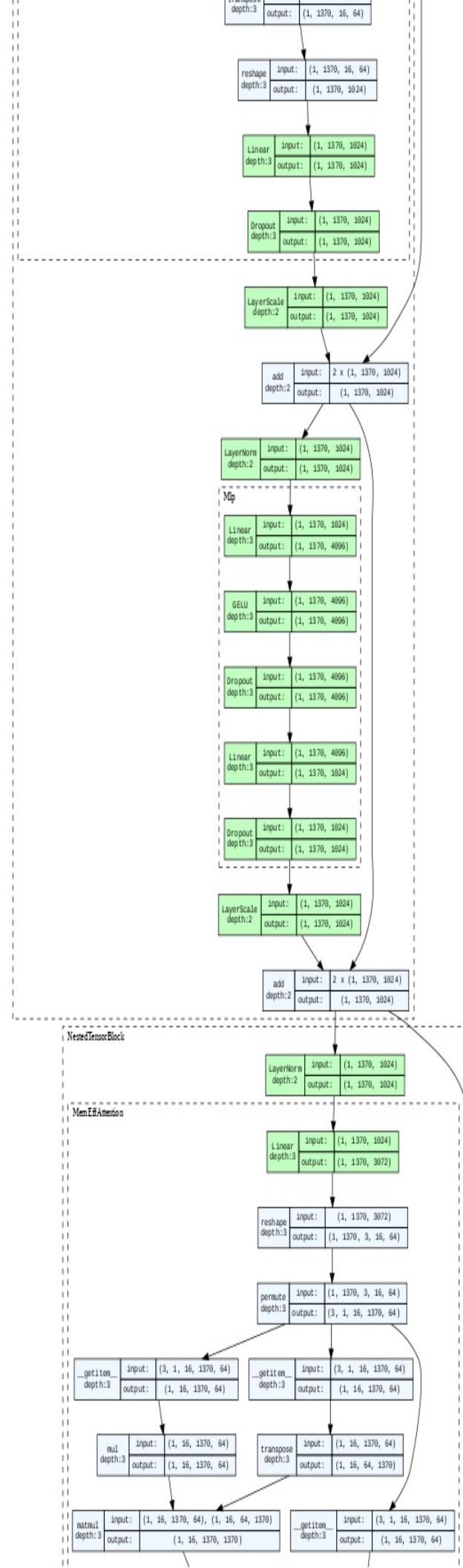












softmax	input:	(1, 16, 1370, 1370)
	output:	(1, 16, 1370, 1370)

Dropout	input:	(1, 16, 1370, 1370)
	output:	(1, 16, 1370, 1370)

matmul	input:	(1, 16, 1370, 1370), (1, 16, 1370, 64)
	output:	(1, 16, 1370, 64)

transpose	input:	(1, 16, 1370, 64)
	output:	(1, 1370, 16, 64)

reshape	input:	(1, 1370, 16, 64)
	output:	(1, 1370, 1024)

Linear	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

Dropout	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

LayerScale	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

add	input:	2 x (1, 1370, 1024)
	output:	(1, 1370, 1024)

LayerNorm	input:	(1, 1370, 1024)
	output:	(1, 1370, 1024)

Mp

Linear	input:	(1, 1370, 1024)
	output:	(1, 1370, 4096)

gelu	input:	(1, 1370, 4096)
	output:	(1, 1370, 4096)

Dropout	input:	(1, 1370, 4096)
	output:	(1, 1370, 4096)

Linear	input:	(1, 1370, 4096)
	output:	(1, 1370, 3072)

Dropout	input:	(1, 1370, 3072)
	output:	(1, 1370, 3072)

LayerScale	input:	(1, 1370, 3072)
	output:	(1, 1370, 3072)

add	input:	2 x (1, 1370, 3072)
	output:	(1, 1370, 3072)

NextTensorBlock

LayerNorm	input:	(1, 1370, 3072)
	output:	(1, 1370, 3072)

Linear	input:	(1, 1370, 3072)
	output:	(1, 1370, 3672)

reshape	input:	(1, 1370, 3672)
	output:	(1, 1370, 3, 16, 64)

permute	input:	(1, 1370, 3, 16, 64)
	output:	(1, 1370, 3672)

