input_ids：词元化的句子，每个词元用id表示
B是batch_size，L是句子长度
input_ids作为输入送入embedding层
> self.embedding = nn.Embedding(vocab_size, d_model, **factory_kwargs)
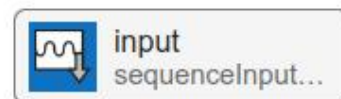vocab_size是词汇表的大小，d_model表示状态空间维数（词向量维数）
接下来送入若干个Block（ResidualBlock, named from mamba-minimal）
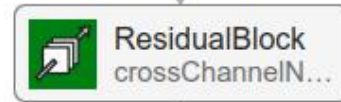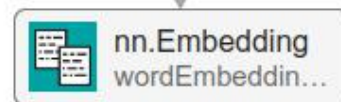得到的结果作归一化（RMSNorm or nn.LayerNorm）后送入线性层得到logits
logits可以用于classify或者decode之后用于generate

```python
1  class MixerModel(nn.Module):
2      def __init__(...) -> None:
3          factory_kwargs = {"device": device, "dtype": dtype}
4          super().__init__()
5          self.residual_in_fp32 = residual_in_fp32
6
7          self.embedding = nn.Embedding(vocab_size, d_model, **factory_kwargs)
8
9          self.fused_add_norm = fused_add_norm
10         if self.fused_add_norm:
11             if layer_norm_fn is None or rms_norm_fn is None:
12                 raise ImportError("Failed to import Triton LayerNorm / RMSNorm kernels")
13
14         self.layers = nn.ModuleList( [ create_block(...) for i in range(n_layer) ] )
15
16         self.norm_f = (nn.LayerNorm if not rms_norm else RMSNorm)(
17             d_model, eps=norm_epsilon, **factory_kwargs
18         )
19
20         self.apply(
21             partial(
22                 _init_weights,
23                 n_layer=n_layer,
24                 **(initializer_cfg if initializer_cfg is not None else {}),
25             )
26         )
```

(B,L)    input   sequenceInput...

(B,L,d_model)    nn.Embedding   wordEmbeddin...

ResidualBlock   crossChannelN...

ResidualBlock   crossChannelN...

ResidualBlock   crossChannelN...

(B,L,d_model)    RMSNorm   layerNormalizat...

(B,L,vocab_size)    nn.Linear   fullyConnected...

d_inner是内部维度，d_inner=d_model * expand（expand=2 by default）
在MambaBlock内部，输入经过第一个线性层之后等分为d_inner的两份，
分别交付两条计算路线。F.silu是激活函数。

**inside MambaBlock**

```python
1  class Mamba(nn.Module):
2      def __init__(...):
3          factory_kwargs = {"device": device, "dtype": dtype}
4          super().__init__()
5          self.d_model = d_model
6          self.d_state = d_state
7          self.d_conv = d_conv
8          self.expand = expand
9          self.d_inner = int(self.expand * self.d_model)
10         ...
11         self.use_fast_path = use_fast_path
12         self.layer_idx = layer_idx
13
14         self.in_proj = nn.Linear(
15             self.d_model,
16             self.d_inner * 2,
17             bias=bias,
18             **factory_kwargs
19         )
20
21         self.conv1d = nn.Conv1d(
22             in_channels=self.d_inner,
23             out_channels=self.d_inner,
24             bias=conv_bias,
25             kernel_size=d_conv,
26             groups=self.d_inner,
27             padding=d_conv - 1,
28             **factory_kwargs,
29         )
30
31         self.activation = "silu"
32         self.act = nn.SiLU()
33
34         ...
35
36         self.out_proj = nn.Linear(
37             self.d_inner,
38             self.d_model,
39             bias=bias,
40             **factory_kwargs
41         )
```

**inside ResidualBlock**

Residual
inputLayer

Input
inputLayer

addition
additionLayer

RMSNorm
layerNormalizat...

MambaBlock
crossChannelN...

Output
crossChannelN...

Residual
crossChannelN...

(B,L,d_model)
Input
inputLayer

(B,L,2*d_inner)
nn.Linear
fullyConnected...

conv1d
convolution1dL...

(B,L,d_inner)
F.silu
functionLayer

(B,L,d_inner)
SSMBlock
crossChannelN...

F.silu
functionLayer

(B,L,d_inner)
multiplication
multiplicationLa...

(B,L,d_model)
nn.Linear
fullyConnected...