

Transformers

The Decoder

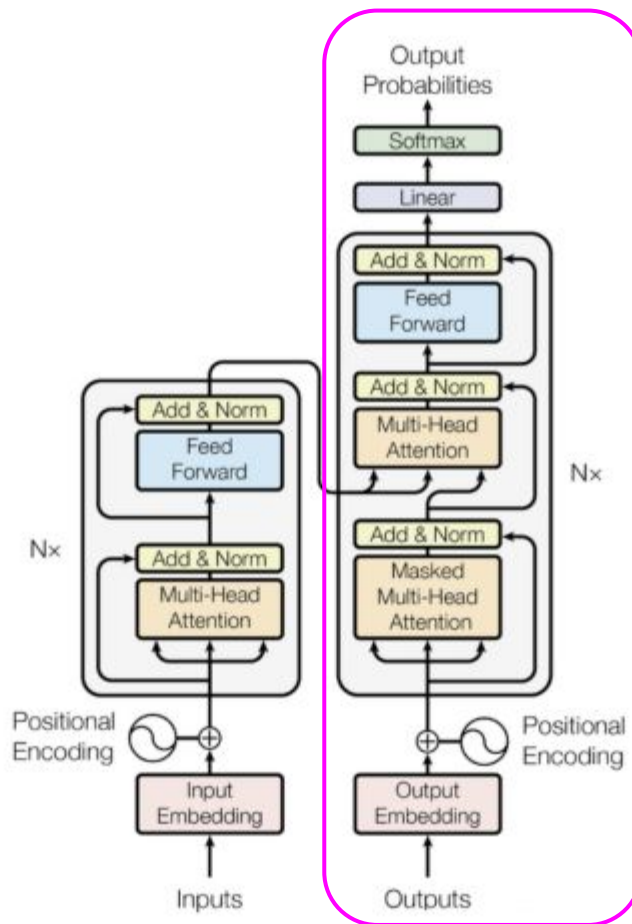
Objectives

- understand the main components of the Transformer
- understand the importance of Attention mechanism
- be able to implement a small/medium Transformer model
- be able to train for long periods of time (several hours)
- be able to do checkpointing and continuous evaluation

GPT

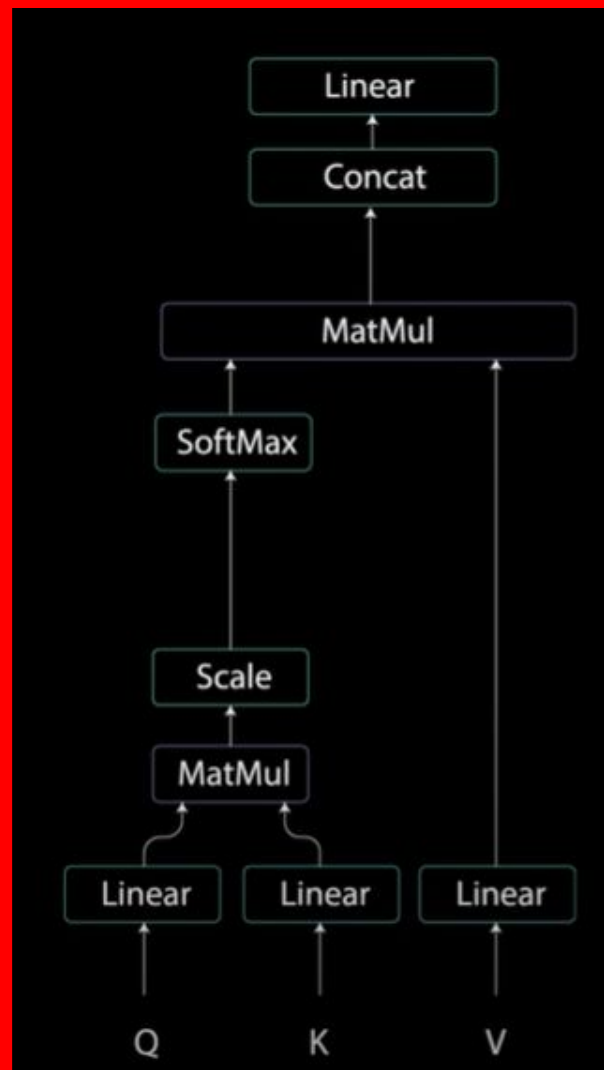
BERT

Encoder



GPT

Decoder



```
class Transformer(torch.nn.Module):
    def __init__(self, vocab_size, input_embedding_dimensions, max_sequence_length):
        super(Transformer, self).__init__()
        #positional encoding
        self.input_embeddings = nn.Embedding(vocab_size, input_embedding_dimensions)
        self.positional_encoding = positional_encoding(max_sequence_length, input_embedding_dimensions)

        #MHA
        self.mha_first_linear_layer = nn.Linear(input_embedding_dimensions, input_embedding_dimensions)
        self.mha_final_linear_layer = nn.Linear(input_embedding_dimensions, input_embedding_dimensions)

        #ADD & NORM LAYER
        self.mha_norm_layer = nn.LayerNorm(input_embedding_dimensions) #

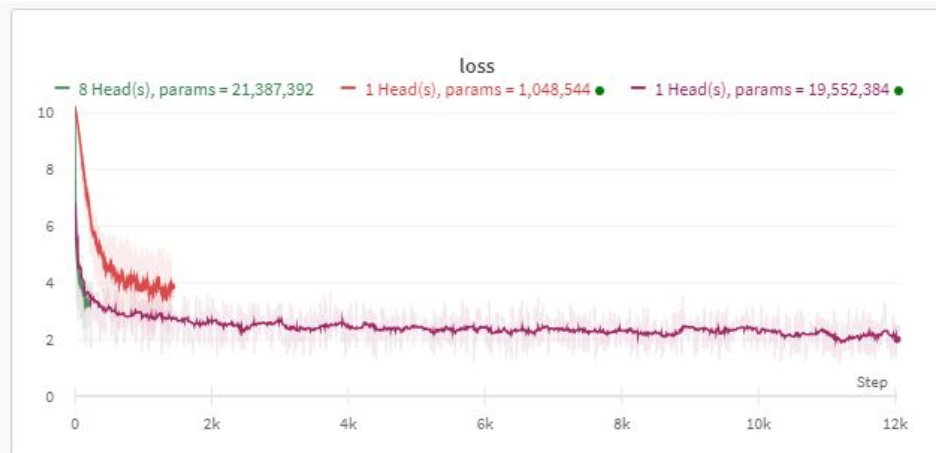
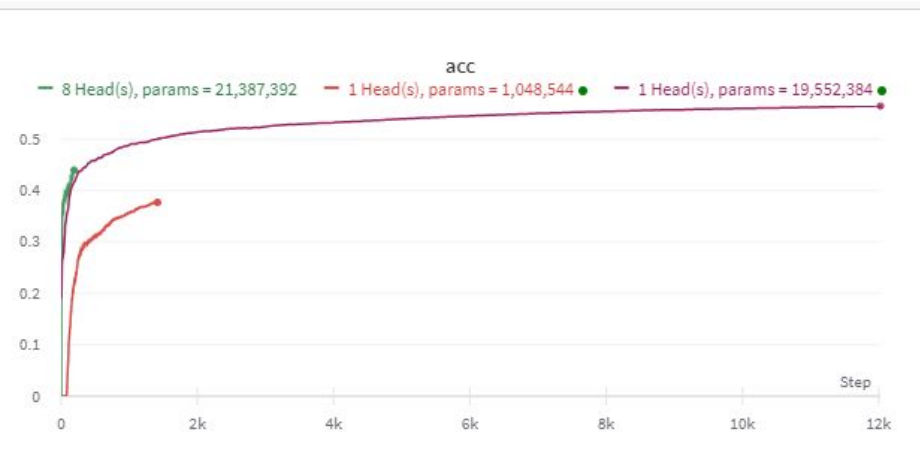
        #FEED FORWARD
        self.ff_first_linear = nn.Linear(input_embedding_dimensions, input_embedding_dimensions)
        self.ff_relu = nn.ReLU()
        self.ff_second_linear = nn.Linear(input_embedding_dimensions, input_embedding_dimensions)
        self.ff_norm_layer = nn.LayerNorm(input_embedding_dimensions) #

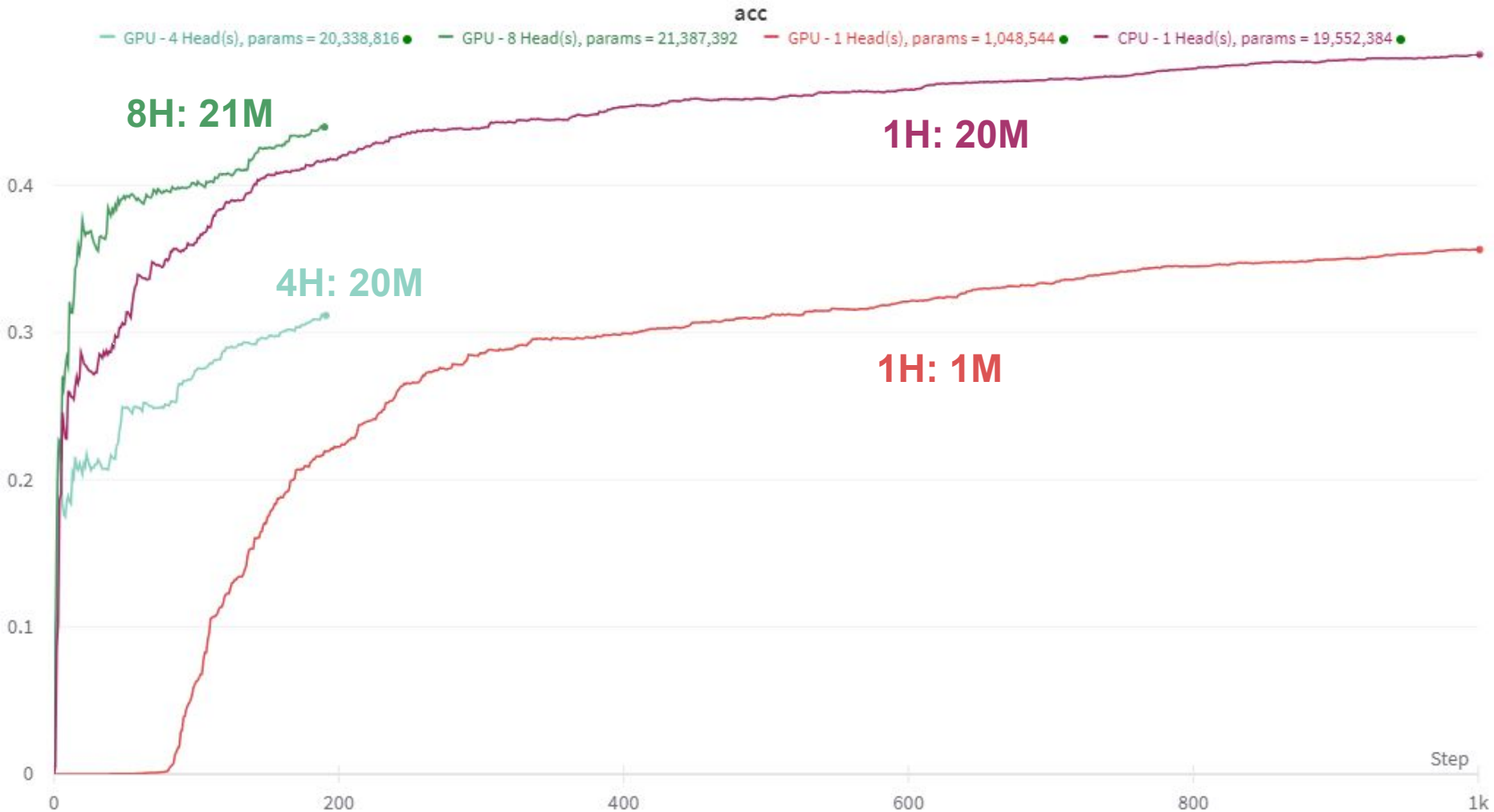
        #final linear layer
        self.final_linear_layer = nn.Linear(input_embedding_dimensions, vocab_size)

    def forward(self, inputs):
        #positional encoding
        input_embeddings = self.input_embeddings(inputs)
        positional_embeddings = input_embeddings + self.positional_encoding

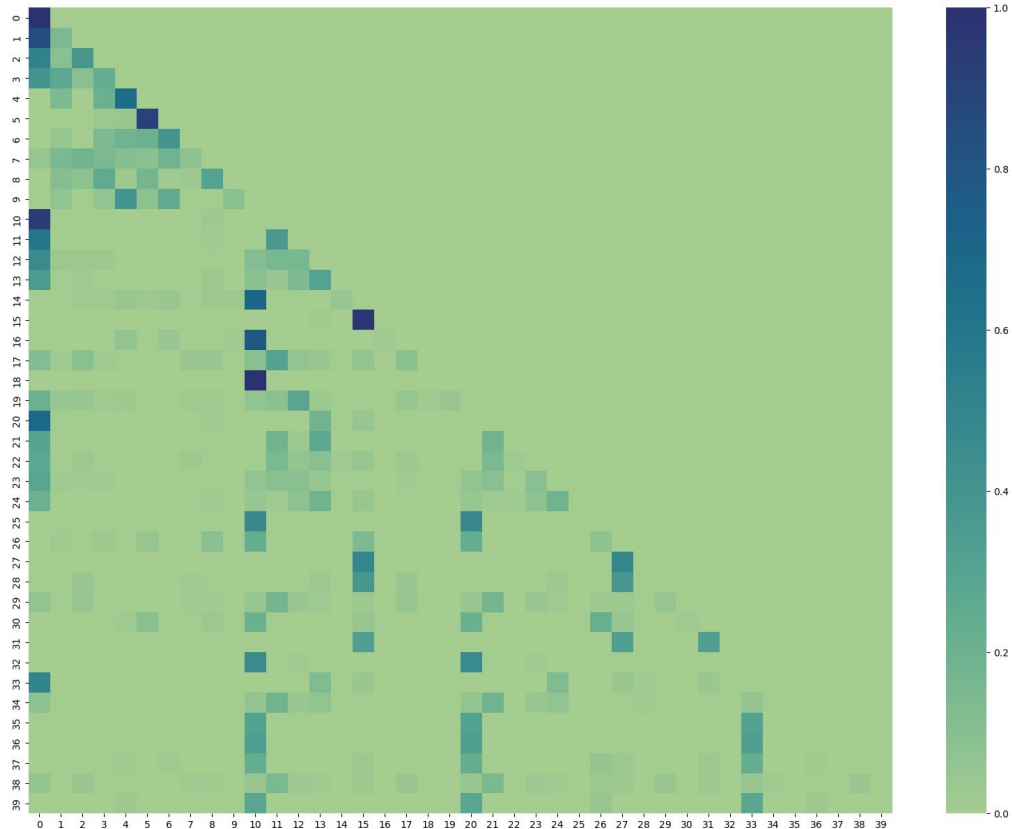
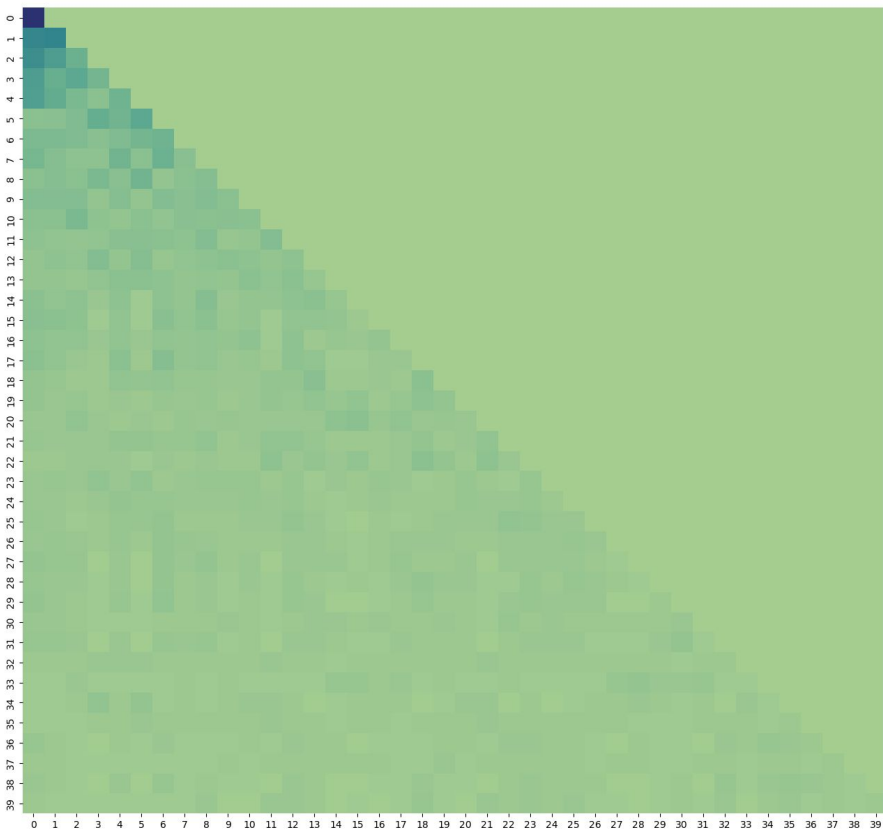
        #MHA
        q, k, v = self.mha_first_linear_layer(positional_embeddings)
        score_matrix = q @ k.transpose(-1, -2)
        score_matrix = score_matrix / (input_embedding_dimensions ** 0.5)
        mask = torch.tril(torch.ones(score_matrix.size(1), score_matrix.size(1)), diagonal=0)
        score_matrix = score_matrix.masked_fill(mask, float('-inf'))
        att_weights = torch.softmax(score_matrix, dim=-1)
        output = att_weights @ v
        final_mha_output = self.mha_final_linear_layer(output)
```

Training

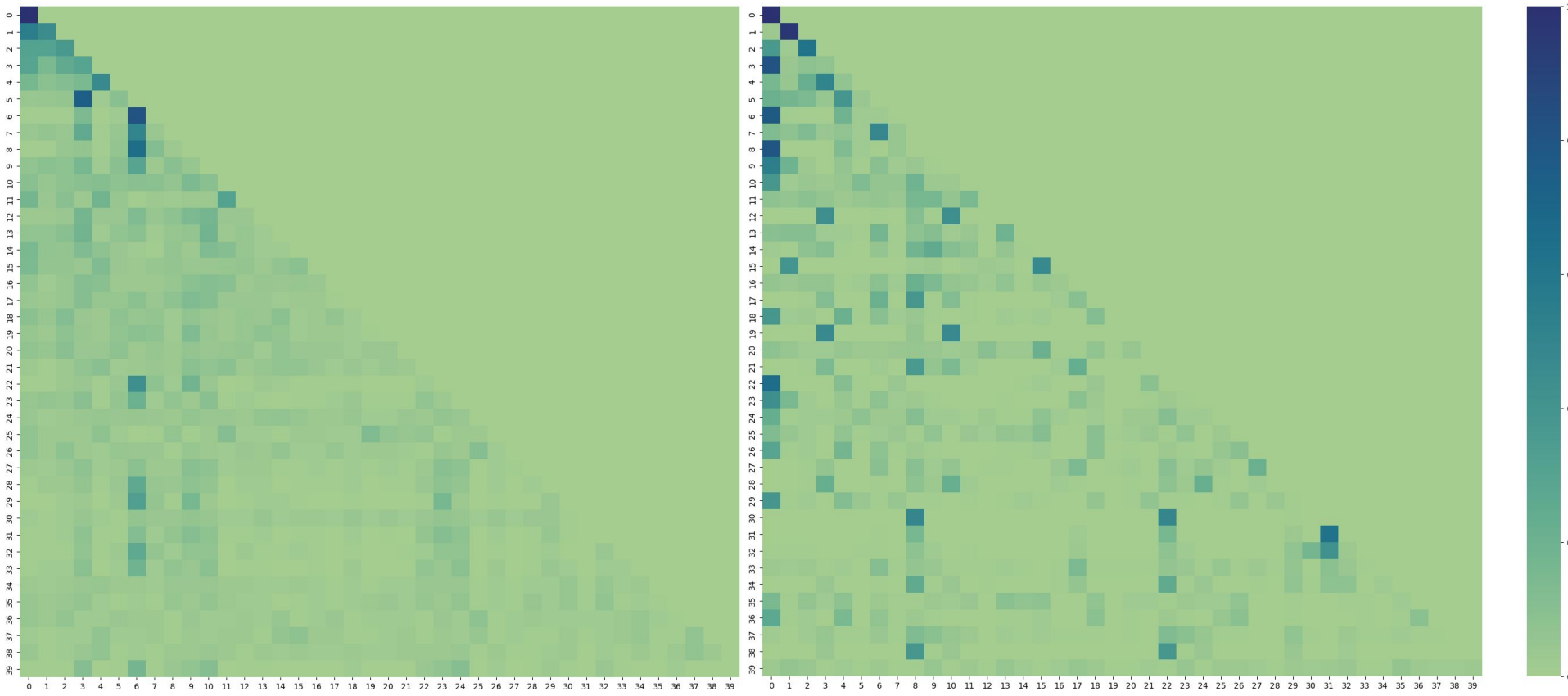




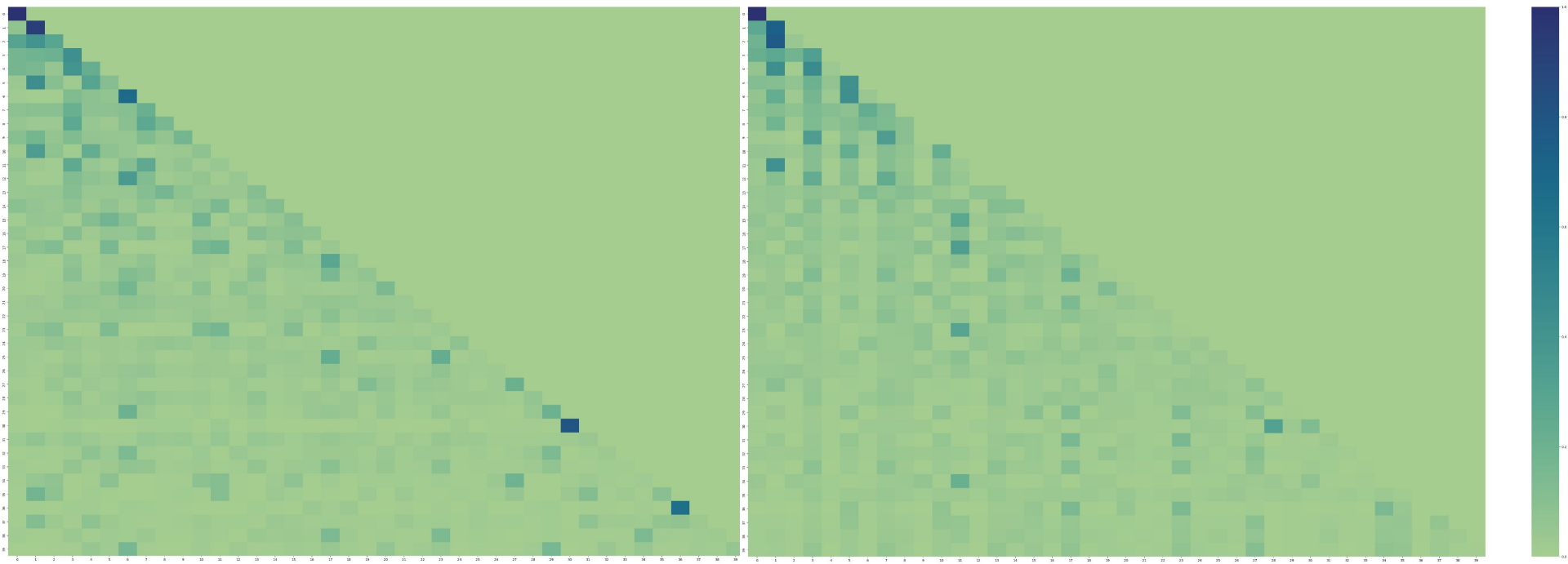
Attention: Initial VS Trained (280K)



Attention: With VS Without Positional Encoding (120k)



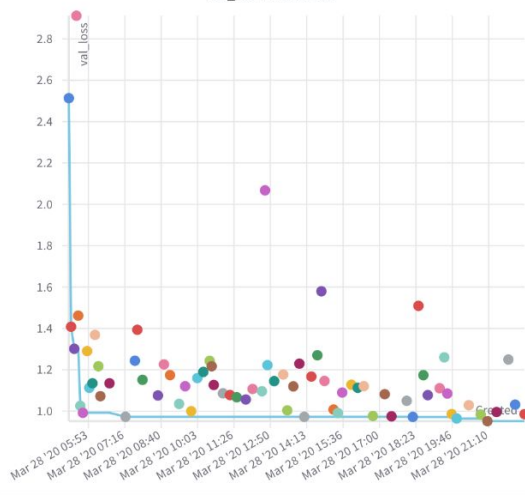
Attention: Different Heads (68,000)



Weights and Biases - Potential

Follow along with a [video tutorial!](#)

val_loss v. created



Parameter importance with respect to

val_loss

Search

Page 1 of 1



Config parameter

Importance

Correlation

features_start

grad_batches

num_layers

batch_size

lr

bilinear

epochs

Runs (18)

Q

11

Name (2 visualized)

● CNN-2, E 1

● CNN-2, E 5

Tables 1

runs_summary["test_predictions"]

Stop

Reset & Automate Columns

Joined On

guess

truth

score_0

score_1

score_2

score_3

score_4

score_5

1-10 of 224

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

GPU

1. Large Model Size:

- A model with a large number of parameters requires more GPU memory and can lead to slower computations.

2. Large Batch Size:

- Larger batch sizes can be beneficial for parallelism but may also lead to increased GPU memory usage and slower computations if the GPU memory is insufficient.

3. Inefficient Memory Usage:

- Inefficient use of GPU memory, such as unnecessary allocations or not releasing memory properly, can lead to slowdowns.

4. Data Transfer Between CPU and GPU:

- Frequent data transfers between CPU and GPU can introduce overhead. Minimizing these transfers can improve performance.

5. Complex Operations:

- Certain complex operation computationally intensive

6. Inefficient Tensor Operations

- Writing inefficient code with redundant computations. For example

7. Insufficient Parallelism:

- If the GPU is not fully utilized, the system may experience bottlenecks in computation. Ensuring that the GPU is fully utilized can improve the performance of parallel processing and reduce the time required to complete the task.

8. Inefficient GPU Kernel Calls:

- Inefficient use of GPU kernel performance. Frameworks implementations for comm

9. Out-of-Memory Errors:

- Running out of GPU memory is essential to manage memory in a model.

[illegible]

What worked in coding the architecture

Building the most simple structure

Moving very slowly

Being a maniac about input/output shapes

Add complexities - One at a time

At every stage, when I think i got it right, I assume i got it wrong

- Talk to others
- Eavesdropping