

РЕ - позиционные эмбединги.

Архитектура трансформера - это архитектура созданная из енкодера и декодера, только их какое то N количество штук. Возьмем классическую схему из хабра,

И на её примере посмотрим работу трансформера. У нас на схеме вводится  $x_1$  и  $x_2$  вектора запросов, далее для каждого из них создается позиционный эмбединг. Что это такое? Позиционный эмбединг - вектора, который хранит положение слова в предложении, он необходим, так как без него будет проблема о незнании положения

$$\epsilon_{ij} = \frac{x_i W^Q (x_j W^K)^T + x_i W^Q (p_{ij}^K)^T}{\sqrt{d}}$$

слова в предложении.

ПЕ разделяются в МХСА на два типа:

1) Абсолютный

2) Относительный

В абсолютном каждый входящий токен на позиции  $i$  сопоставляется с тренировочным ембединг вектором, который показывает строку в матрице R с размером [tokens, dims].

R - тренировочная матрица, инициализированная значениями (0,1)

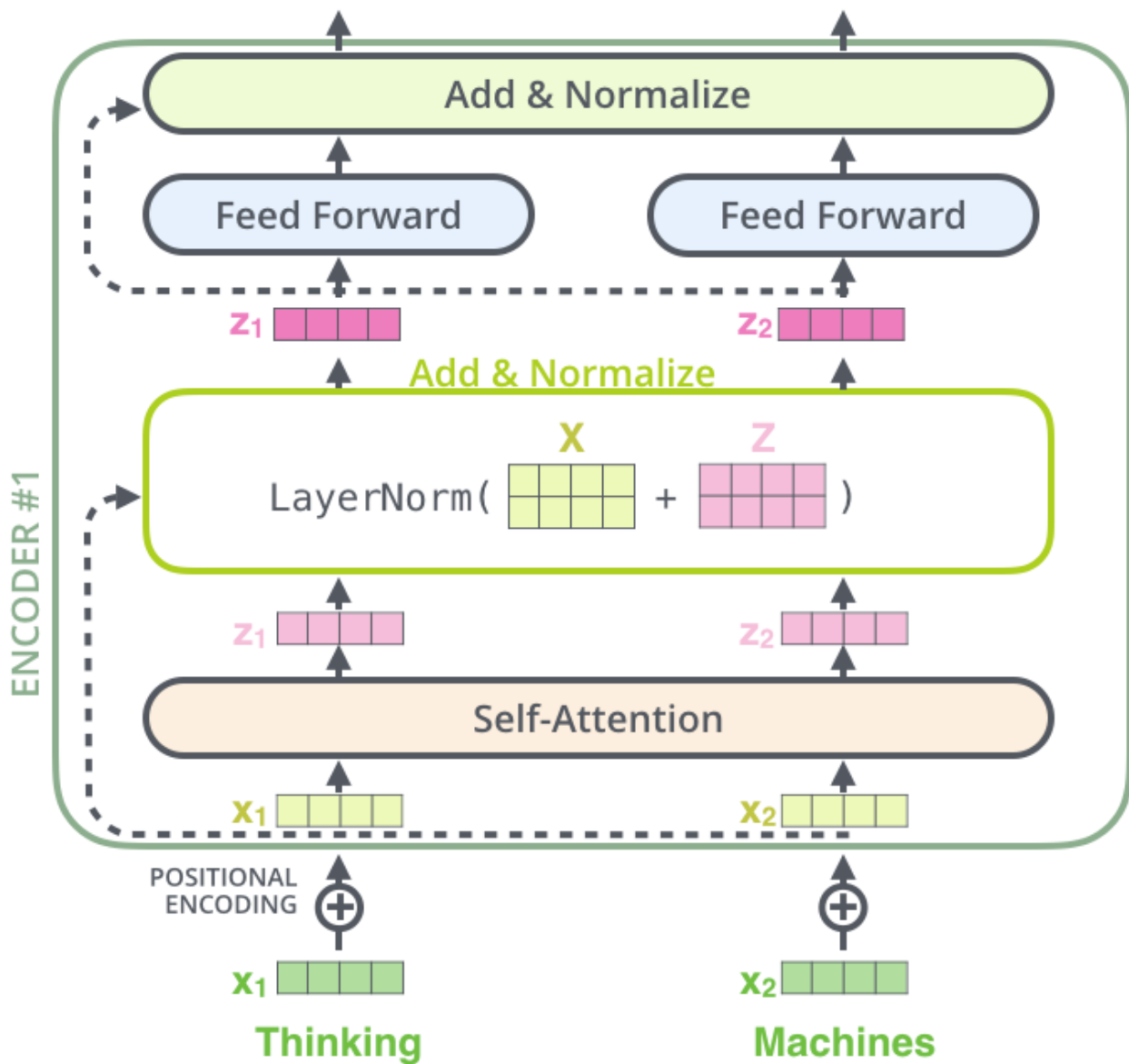
$$att = softmax(\frac{1}{\sqrt{dim}}(QK^T + QR))$$

Относительный ПЕ пресдавляет собой дистанцию между токенами. Но в них могут оказаться трудности с размерностями матрицами, там может быть маленькое

$$\epsilon_{ij} = \frac{x_i W^Q (x_j W^K)^T + x_i W^Q (p_{ij}^K)^T}{\sqrt{d}}$$

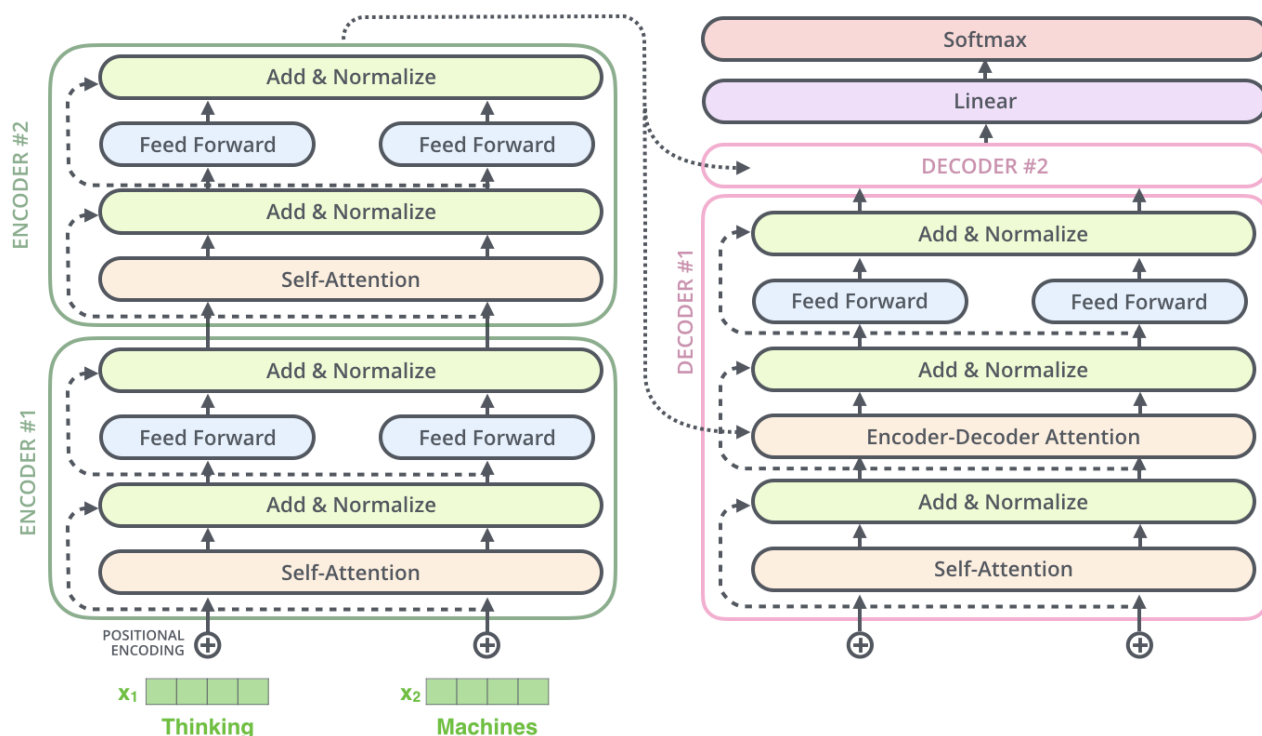
несоответствие с матрицами.

Обычно используют функции косинуса и синуса в функции ПЕ, почему так углубляться не будем.



Получается мы рассматриваем энкодер, в энкодер вводятся векторы  $x_1$  и  $x_2$ , далее они прогоняются через слой внутреннего внимания добавляются к прошлым векторам и нормализуются, далее они проходят через обычную сеть прямого распространения и также суммируются и нормализуются.

## Декодер



Энкодер и декодер были рассмотрены ранее, но их полновязная работа в трансформерах представлена на этой схеме.

Линейный слой – это простая полновязная нейронная сеть, которая переводит вектор, созданный стеком декодеров, в значительно больший вектор, называемый логит вектором (logits vector).

Пусть наша модель знает 10 тысяч уникальных английских слов («выходной словарь» нашей модели), которые она узнала из обучающего корпуса. Это означает, что наш логит вектор будет иметь 10 000 ячеек в ширину – каждая ячейка соответствует коэффициенту одного уникального слова. Таким образом мы интерпретируем выход нашей модели с помощью линейного слоя.

Слой софтмакс переводит этот показатель в вероятности (положительные числа, сумма которых равна 1). Выбирается ячейка с наиболее высокой вероятностью и на выход данного временного отрезка подается соответствующее слово.