

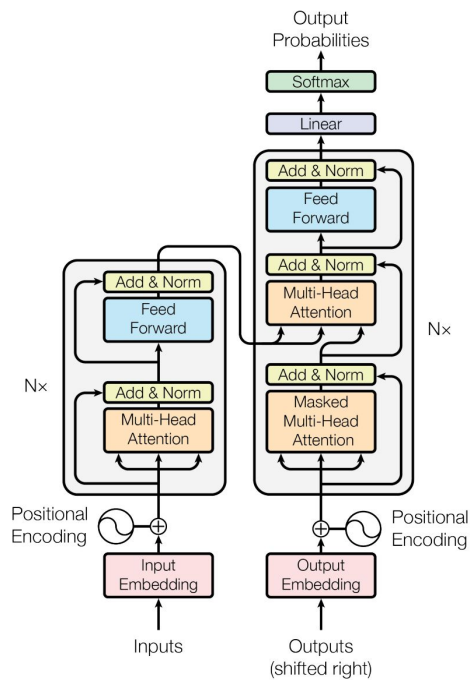
Анализ временных рядов

Азиз Темирханов

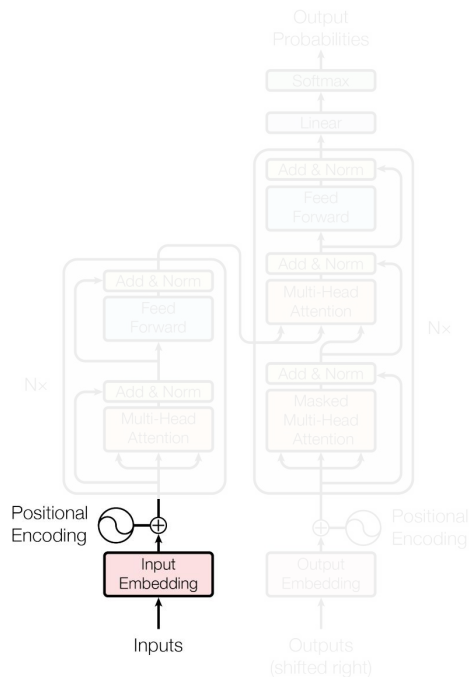
План лекции

- Вспомним механизм внимания
- Узнаем про архитектуры с трансформером
- Поговорим о фундаментах

Transformer



Transformer



Входной текст может быть символами, словами, «токенами»:

"The detective investigated" -> [The_] [detective_] [invest] [igat] [ed_]

Токены — это индексы в «словаре»:

[The_] [detective_] [invest] [igat] [ed_] -> [3 721 68 1337 42]

Каждая запись словаря соответствует выученному вектору размерности d_{model} .

[3 721 68 1337 42] -> [[0.123, -5.234, ...], [...], [...], [...], [...]]

Позиционный энкодинг

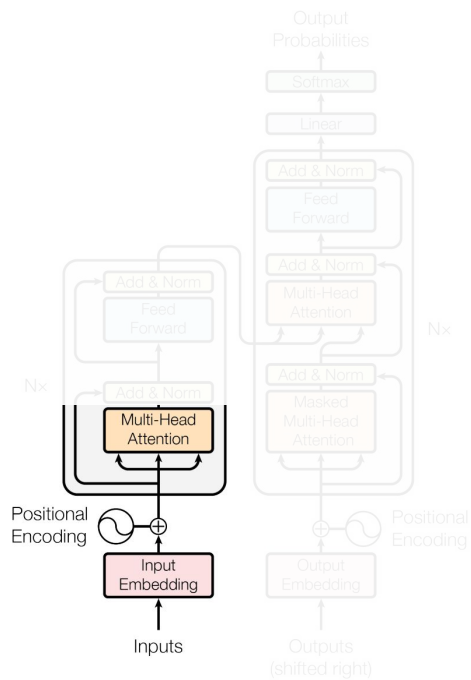
Помните: attention инвариантен к перестановкам, а язык — нет!

("The mouse ate the cat" vs "The cat ate the mouse")

Нужно закодировать позицию каждого слова; просто добавьте что-нибудь.

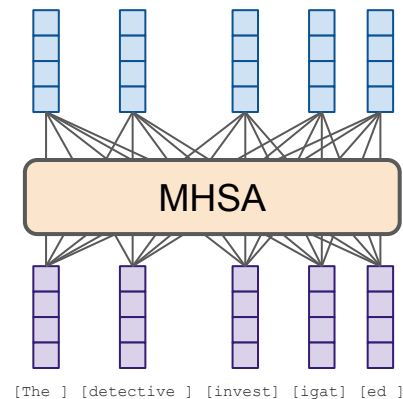
Думайте так: [The_] + 10 [detective_] + 20 [invest] + 30 ... но умнее.

Transformer



То есть входная последовательность используется для построения **queries**, **keys** и **values**!

Каждый токен может «оглядываться» на весь вход и решать, как обновить своё представление, исходя из того, что он видит.



Attention

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \end{matrix}$$
$$= \begin{matrix} \text{Z} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \end{matrix}$$

<http://jalammar.github.io/illustrated-transformer>

MLP

Простой MLP применяется к каждому токenu:

$$z_i = W_2 \text{GeLU}(W_1 x + b_1) + b_2$$

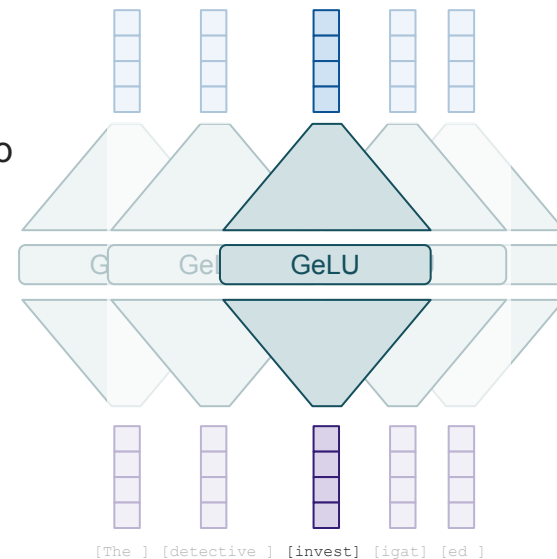
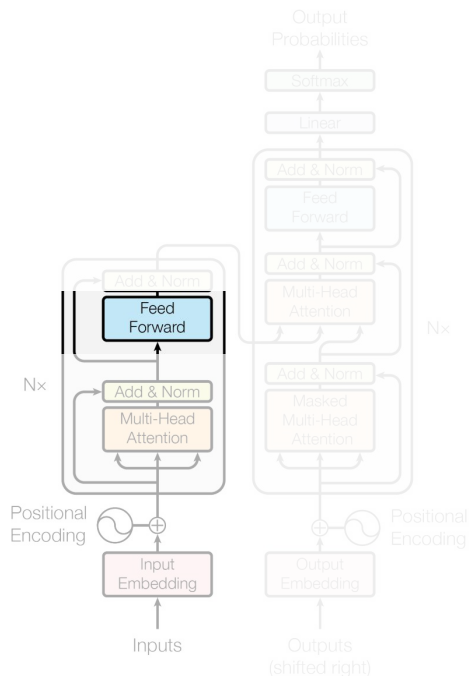
Представьте, что каждый токен самостоятельно размышляет над тем, что он наблюдал ранее.

Есть также некоторые слабые доказательства того, что именно здесь хранится «знание о мире».

В нём содержится основная часть параметров.

Когда люди создают гигантские модели и разреженные/MOE-модели, именно это и становится гигантским.

Некоторые люди любят называть это 1x1 свертками



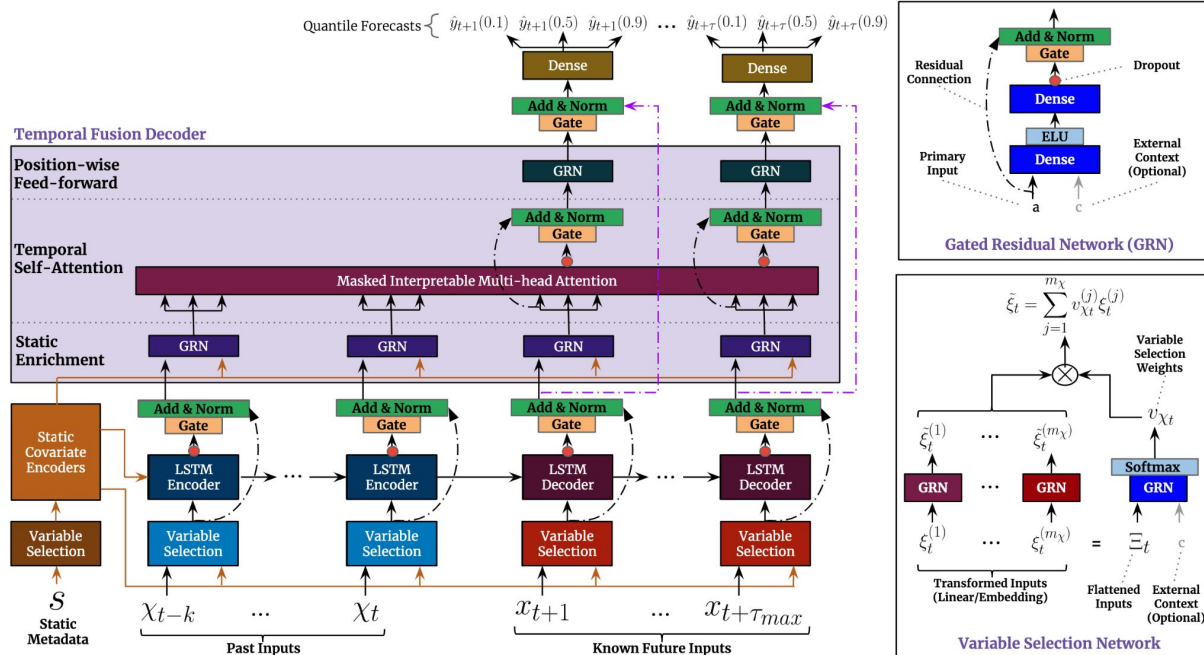
Transformer-модели

Temporal Fusion Transformer (TFT, 2019, 2,148 citations)

Основная идея: многогоризонтное прогнозирование со смешанными ковариатами (статическими, известными на будущее, наблюдаемыми из истории), объединяющее рекуррентные слои (локальные зависимости) + интерпретируемый self-attention + отбор признаков + гейтинг.

Результаты: сообщается о существенных улучшениях по сравнению с бенчмарками на нескольких реальных наборах данных, а также приводятся примеры интерпретируемости и сценарии ее использования.

Ограничения / область применимости: supervised модель

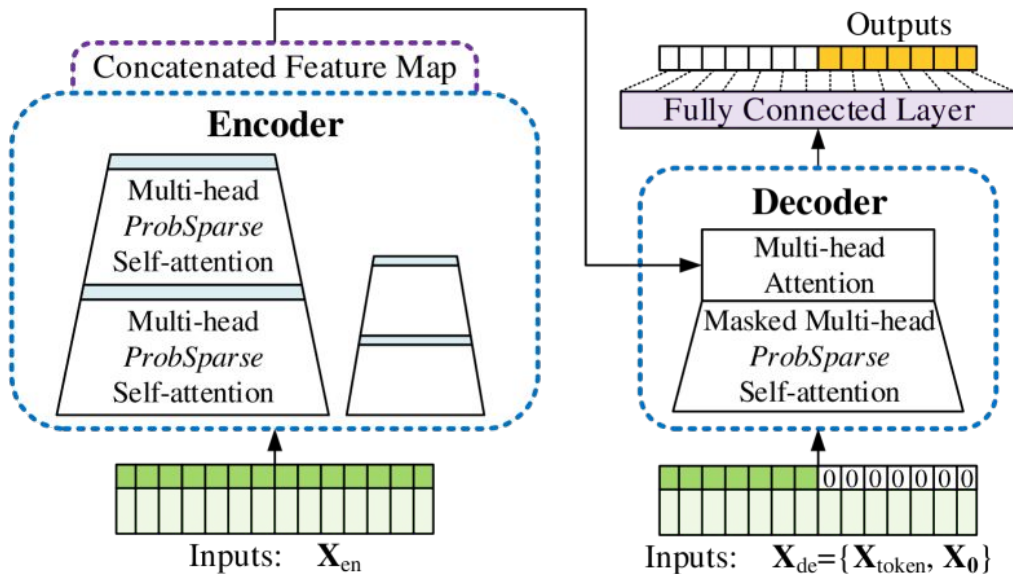


Informer (2020, 6368 citations)

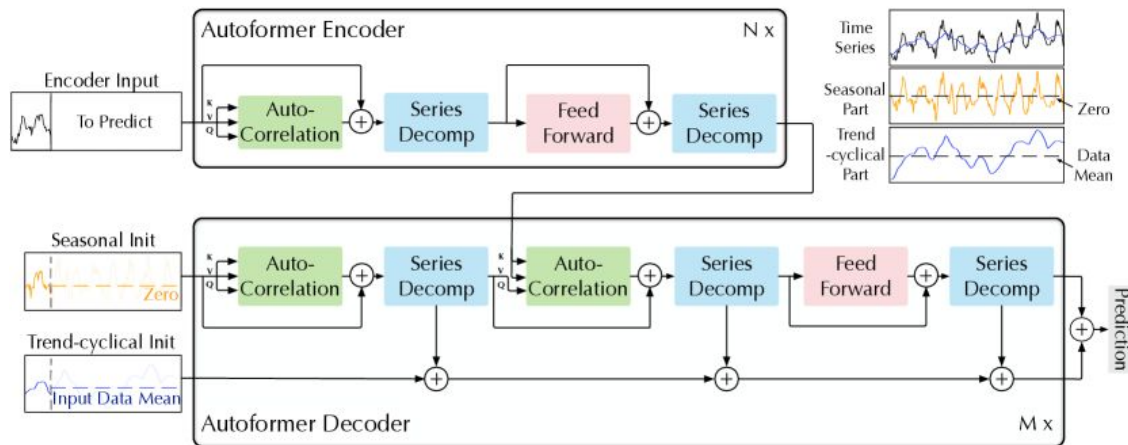
Основная идея: сделать трансформеры практичными для прогнозирования временных рядов по длинным последовательностям (LSTF) с помощью внимания ProbSparse (приблизительно $O(L \log L)$), а также дистилляции внимания, чтобы справляться с чрезвычайно длинными входами.

Результаты: демонстрирует сильные результаты в LSTF при снижении затрат по времени и памяти по сравнению со стандартным вниманием.

Ограничения, отмеченные позже: критика в работе LTSF-Linear / DLinear утверждает, что self-attention, инвариантное к перестановкам, может терять информацию о временном порядке, ставя под сомнение многие варианты трансформеров, используемые в LSTF.



Autoformer (2021, 3726 citations)

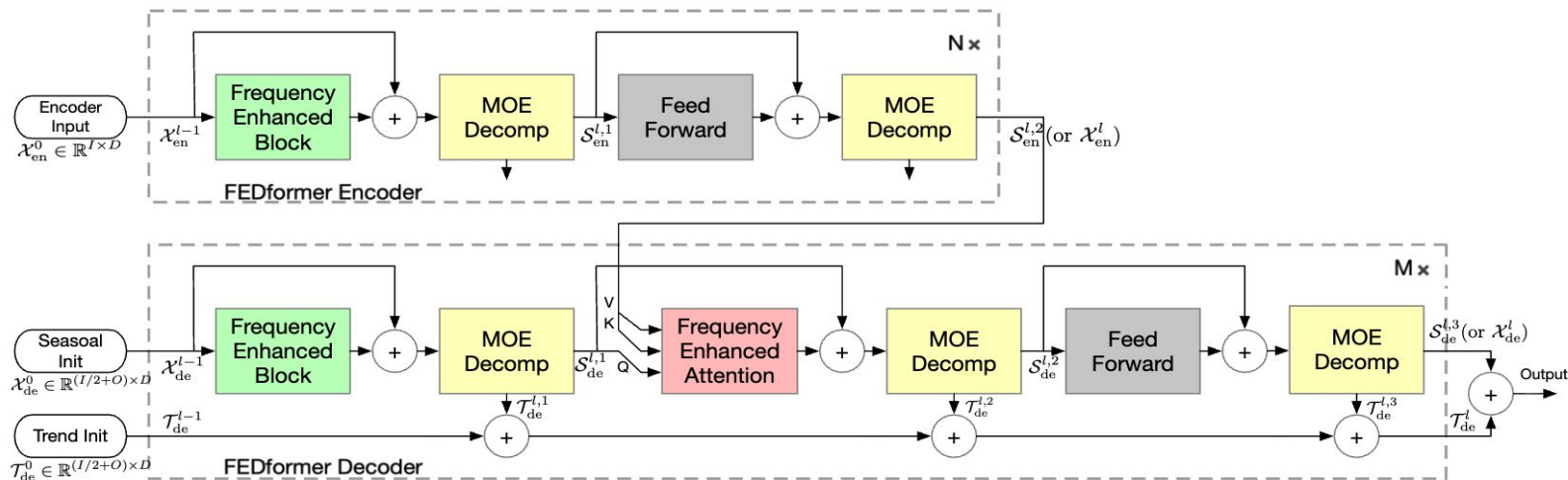


Основная идея: добавить явную декомпозицию сезонность–тренд и заменить self-attention механизмом авто-корреляции, который агрегирует зависимости на уровне подрядов (используя периодичность).

Результаты: заявляет SOTA в долгосрочном прогнозировании с $\sim 38\%$ относительным улучшением на шести бенчмарках.

Ограничения / предположения: ключевой механизм построен вокруг периодичности / структуры автокорреляции ряда — это сильное индуктивное смещение, которое может подходить не всем доменам одинаково хорошо.

FEDformer (2022, 2542 citations)



Основная идея: объединить декомпозицию сезонность–тренд с «частотно-усиленным» трансформером, исходя из того, что многие ряды имеют разреженное представление в частотной области (Фурье / вейвлеты), и тем самым добиться линейной сложности по длине последовательности.

Результаты: сообщает о снижении ошибки примерно на ~14.8% (многомерные ряды) и ~22.6% (одномерные ряды) по сравнению с предыдущим SOTA на шести бенчмарках.

Ограничения: явно опирается на предпосылку «разреженности в частотном базисе» — качество может зависеть от того, насколько это соответствует данным.

DLinear

DLinear (2023, 4319 citations)

Идея: использовать простые линейные модели

$$\hat{X}_i = W X_i \quad W \in \mathbb{R}^{T \times L}$$

Предобработка:

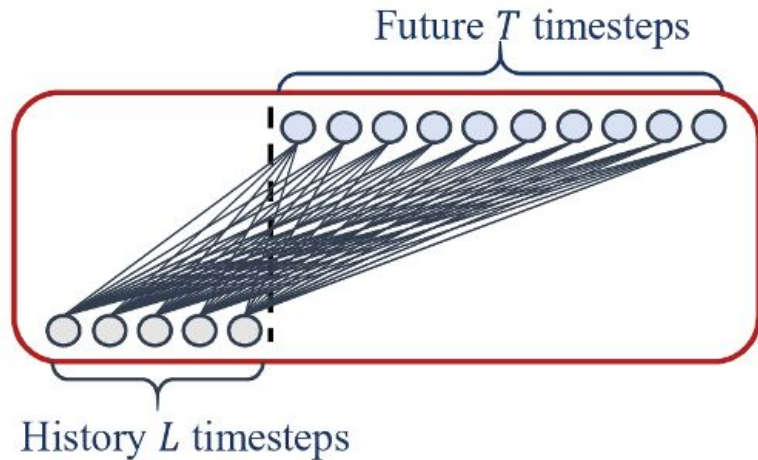
- Декомпозиционный модуль из [Autoformer](#)
- Извлечение тренда скользящим средним

$$\mathcal{X}_t = \text{AvgPool}(\text{Padding}(\mathcal{X}))$$

$$\mathcal{X}_s = \mathcal{X} - \mathcal{X}_t,$$



отдельные модели



Код: <https://github.com/vivva/DLinear>

Статья: <https://arxiv.org/abs/2205.13504>

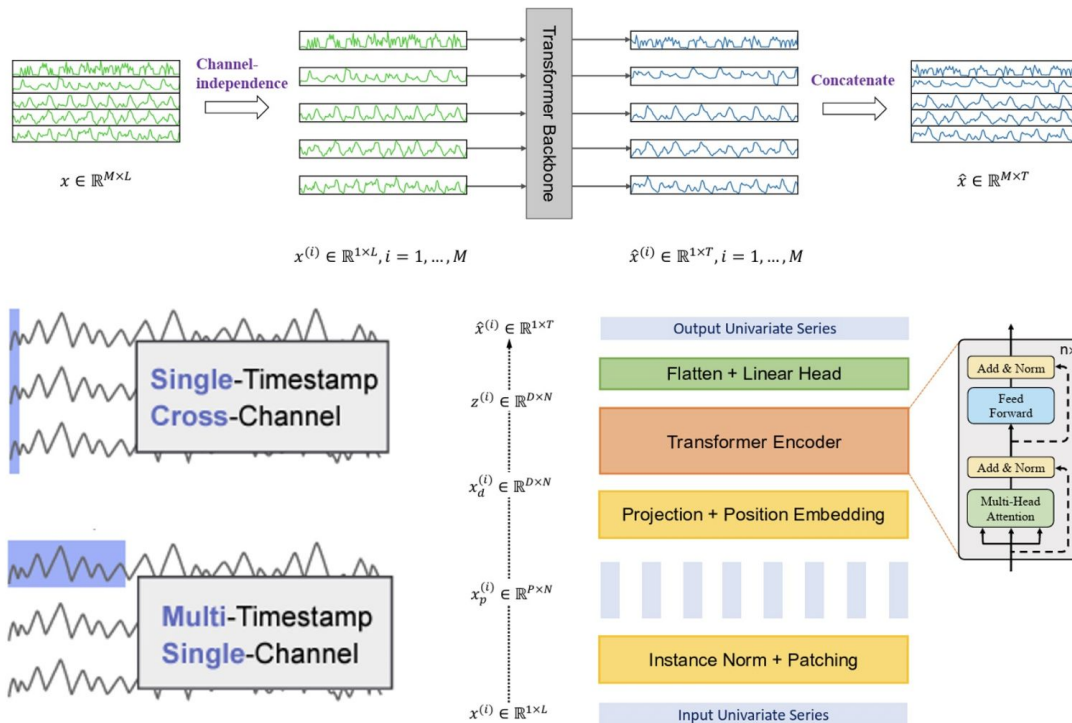
PatchTST (2022, 2673 citations)

Идея: одно наблюдение в конкретный момент времени несет мало информации
(как один пиксель в картинке) - лучше рассматривать “патчи”

Нормализуем и патчим входную последовательность

Получаем эддинги, которые отправляем в трансформерный encoder

Используем ten + линейную голову для получения предсказаний

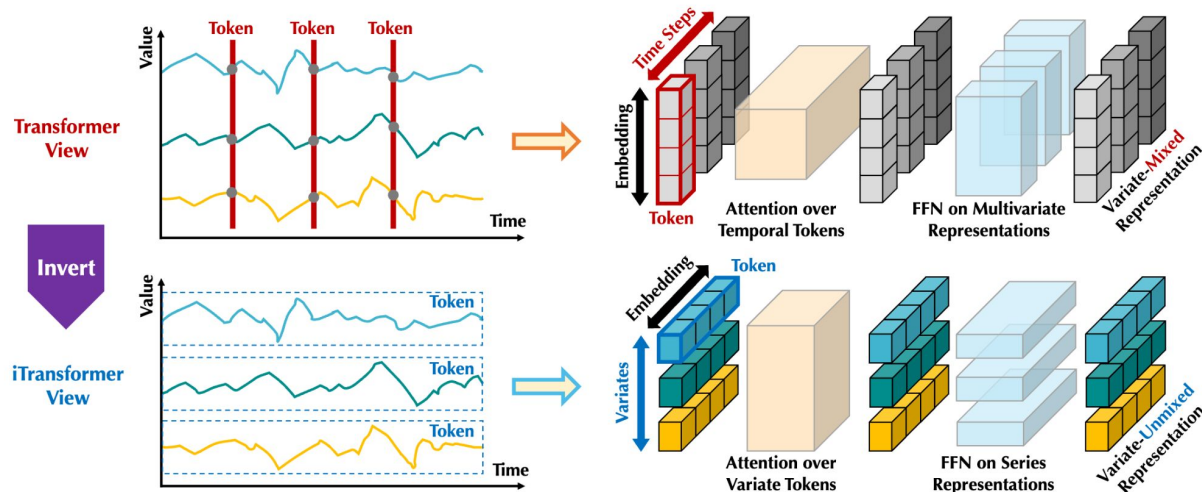


iTransformer (2023, 1241 citations)

Основная идея: инвертировать токенизацию: рассматривать переменные (каналы) как токены и применять внимание *между* каналами, решая проблемы подходов, где временные токены/эмбединги смешивают несколько каналов и плохо масштабируются при длинном окне истории.

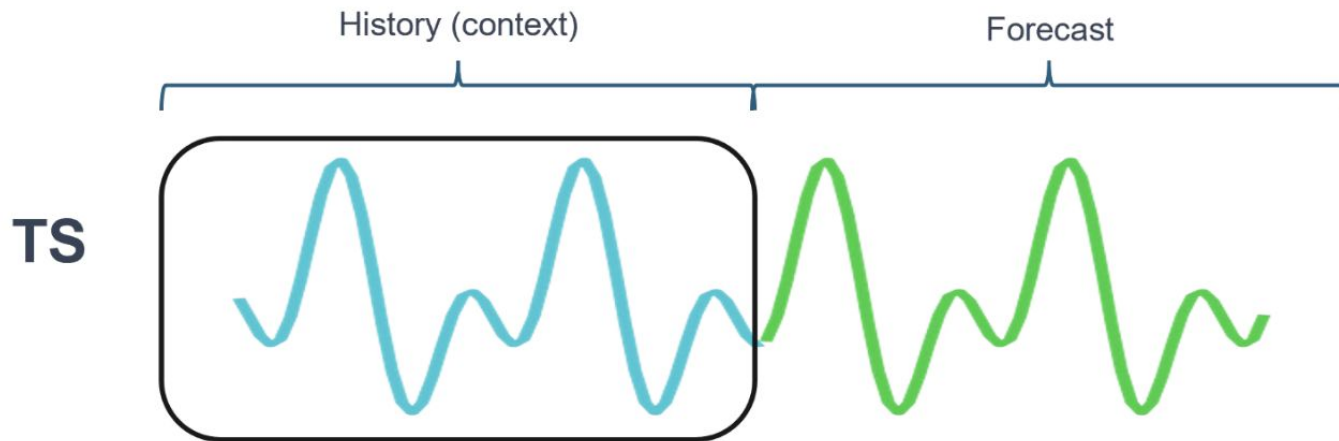
Результаты: заявляет SOTA на реальных датасетах, лучшую обобщаемость между каналами и более хорошую работу с произвольной длиной окна наблюдений.

Зачем: прямой ответ эпохе после DLinear: сохранить стандартные блоки Transformer, но изменить то, *над чем* работает внимание (по каналам, а не по времени).



Foundational Models

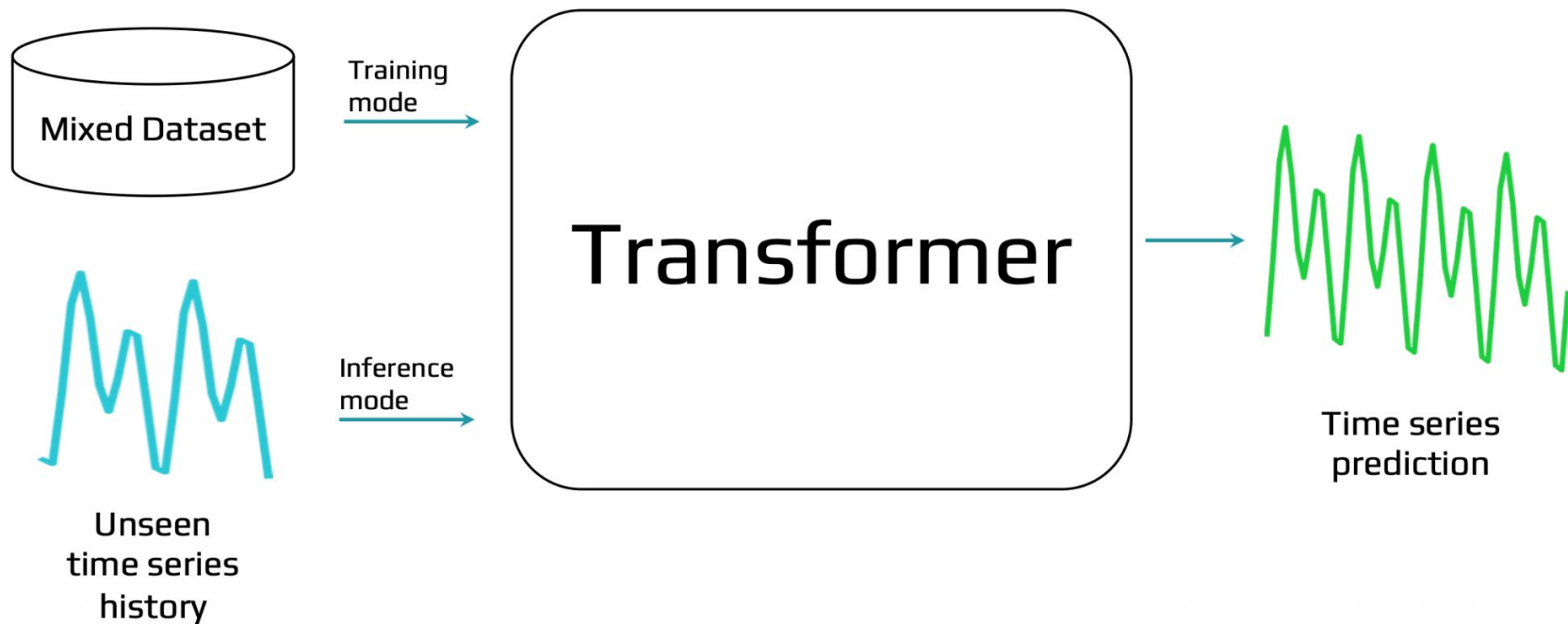
Что такое Zero-Shot?



Zero-shot: Train the model to predict labels for new data **without training on the target dataset**, based on patterns identified in the unrelated data.

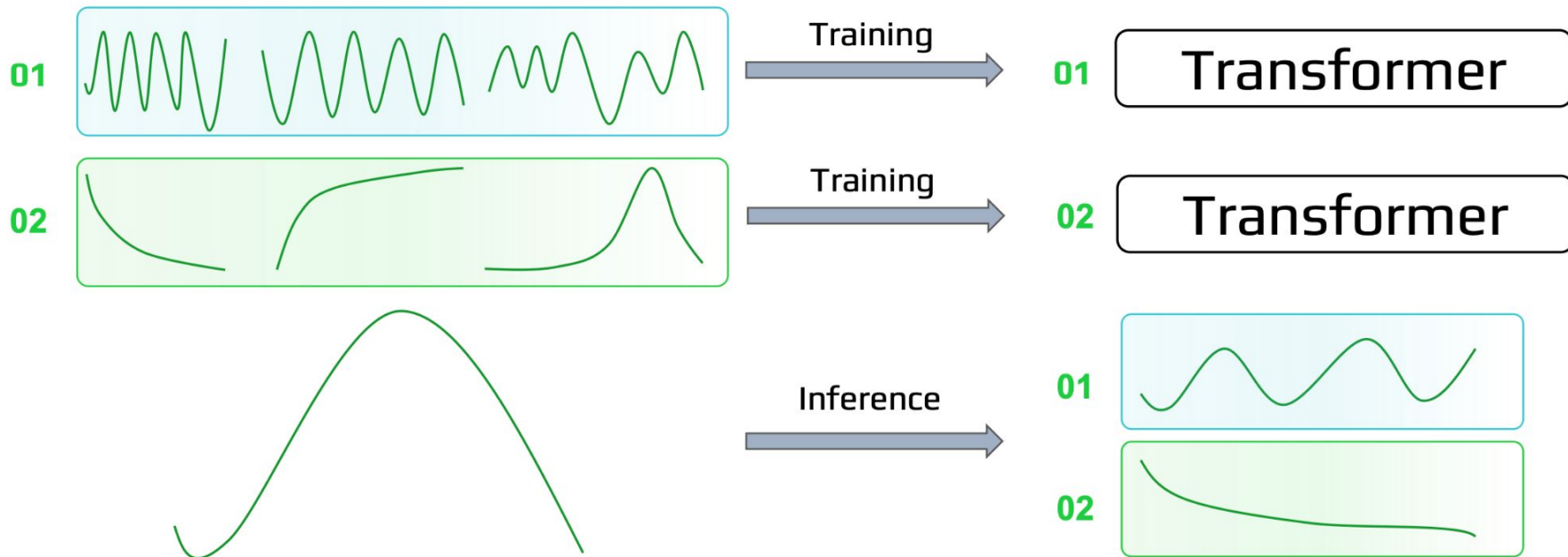
Что такое Zero-Shot?

Zero-shot модели для временных рядов – это в основном трансформеры. Также они требуют большой разнообразный датасет для предобучения.



Что такое Zero-Shot?

Датасеты для обучения



LLMTime (2023, 886 citations)

Идея и Архитектура: кодируем последовательность как текст и используем любую из доступных LLM.

Как кодировать временной ряд?

Rescaling

$$x_t \rightarrow \frac{x_t - b}{a} \quad \left\{ \begin{array}{l} b = \min_t x_t - \beta \cdot (\max_t x_t - \min_t x_t) \\ a = a\text{-percentile}(x_1 - b, x_2 - b, \dots, x_T - b) \end{array} \right.$$

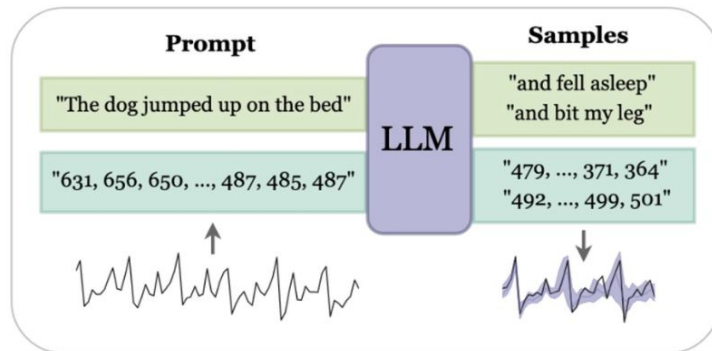
Type-changing 0.123, 1.23, 12.3, 123.0 → "12,123,1230,12300"

Tokenization

"151,167,...,267"

"151,167,...,267"

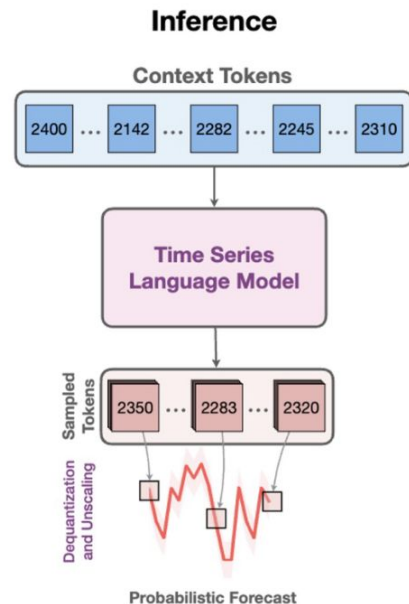
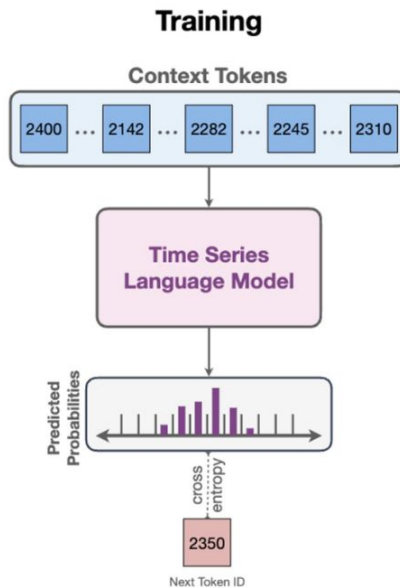
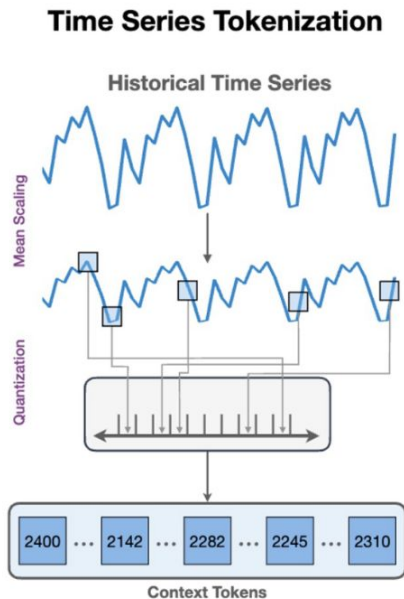
$\alpha=0.95, \beta=0.3$



Упрощенно, задача прогноза — это та же задача продолжения последовательности

CHRONOS (2024, 797)

Идея и Архитектура: кодируем последовательность как текст и используем любую из доступных LLM (T5). Только теперь с дообучением на синтетических и реальных данных с применением TSMixup аугментации (~84 млрд).

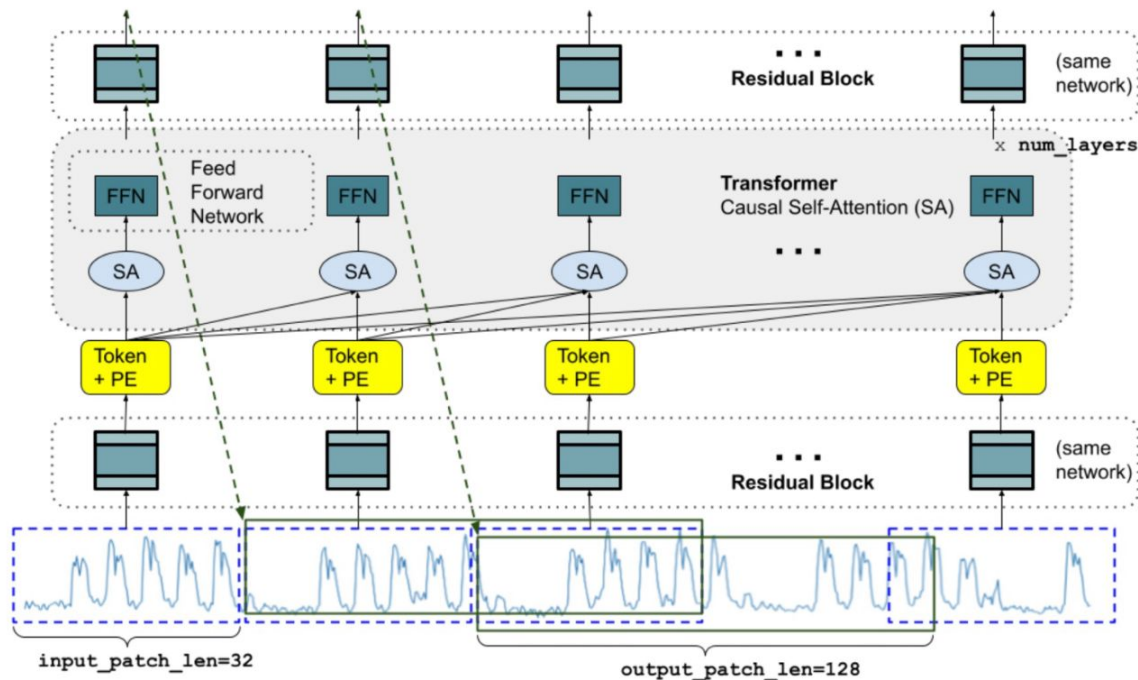


TimesFM (2024, 709)

Идея: обучаем трансформерный decoder на большом количестве синтетических и реальных рядов (~370 млрд).

Задачи: forecasting (+ quantile)

Постановка: univariate
(+ exogenous features)

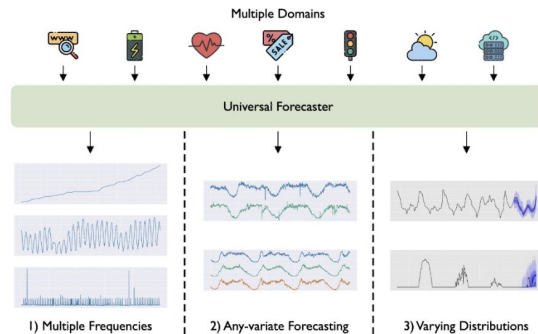


Moirai (2024, 467)

Идея: маскированная encoder-only модель, предобученная на рядах с различной грануляцией (~27 млрд наблюдений).

Задача: probabilistic forecasting

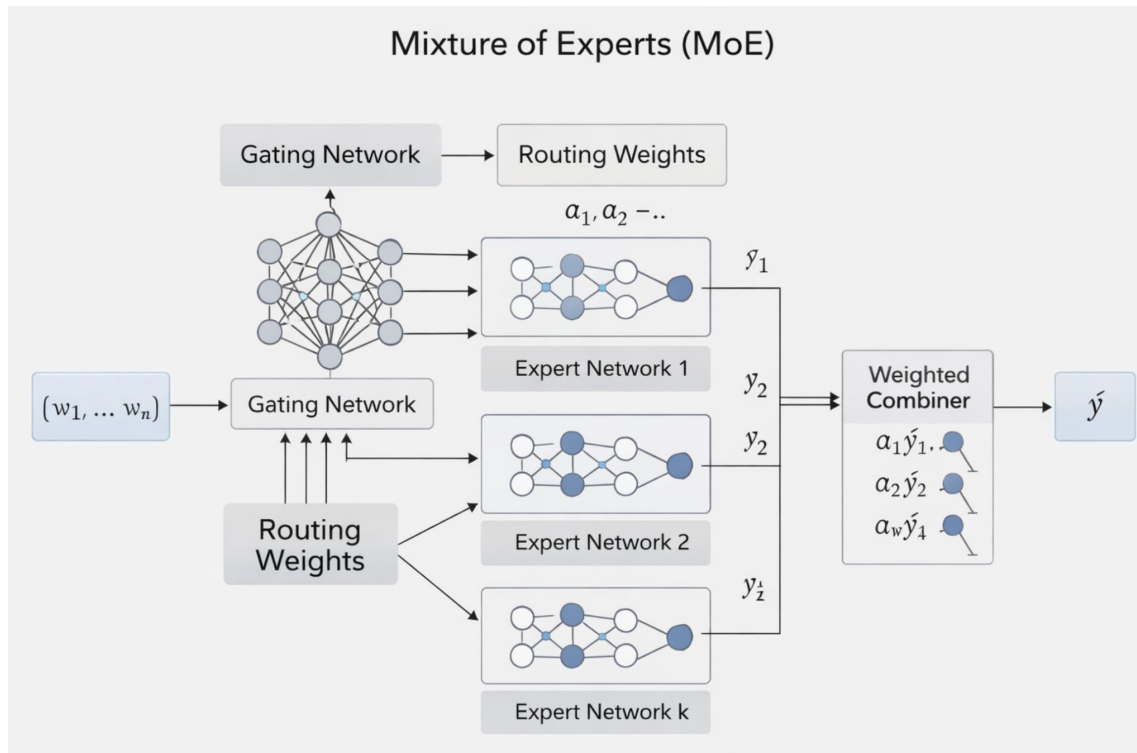
Архитектура: RoPE эмбедингг, flatten для multivariate данных, благодаря чему может обрабатывать данные с произвольной размерностью.



9 доменов. Разная гранулярность — от года до секунды



MoE



Moirai-MOE

Идея: частота — плохой прокси типа паттерна; вместо этого нужна **автоматическая специализация** внутри модели на уровне токенов.

Метод: единая проекция I/O + **sparse**

Mixture-of-Experts внутри Transformer, чтобы разные эксперты ловили разные режимы/паттерны (token-level specialization).

Результат: в экспериментах на десятках датасетов — улучшения над Moirai и конкурентность/превосходство над другими TSFM при меньшем числе активных параметров.

Проблемы: МоЕ сложнее обучать/дебажить (балансировка экспертов, маршрутизация); возможны “мертвые” эксперты и нестабильность при сдвиге домена.

