

Locality-Sensitive Hashing (Reformer)

Підготувала Курченко Лилия, група 191

Проблемы трансформеров

- N слоёв – в N раз больше памяти (хранение активации для обратного распространения)
- Промежуточные слои FF часто достаточно большие.

Проблемы трансформеров

- Матрица attention на последовательностях длины L часто требует $O(L^2)$ как в памяти, так и во времени.

Как их решить?

- Реформер

Что нового в реформере?

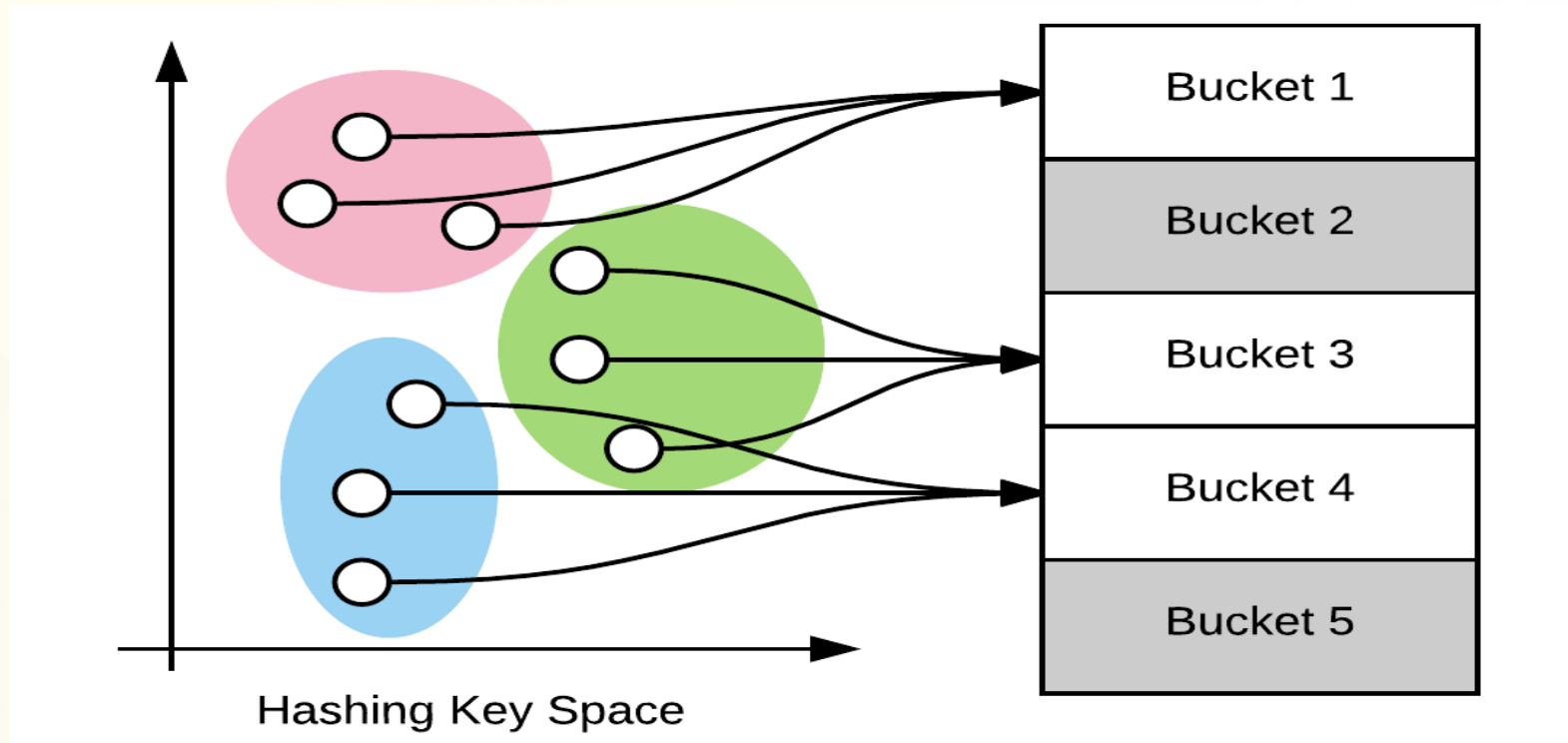
- скалярное произведение attention \rightarrow LSH
 $O(L^2) \rightarrow O(L \log L)$
- Стандартные остаточные блоки \rightarrow
обратимые остаточные слои
N активаций \rightarrow 1 активация сохраняется

Locality-Sensitive Hashing Attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V}$$

- Нас интересуют только самые большие элементы в $\mathbf{Q}\mathbf{K}$.
- Для $q_i \in \mathbf{Q}$ ищем ближайшие к q_i строки в \mathbf{K} .
- Для быстрого поиска используем Locality-Sensitive Hashing (LSH) в attention

Locality-Sensitive Hashing Attention

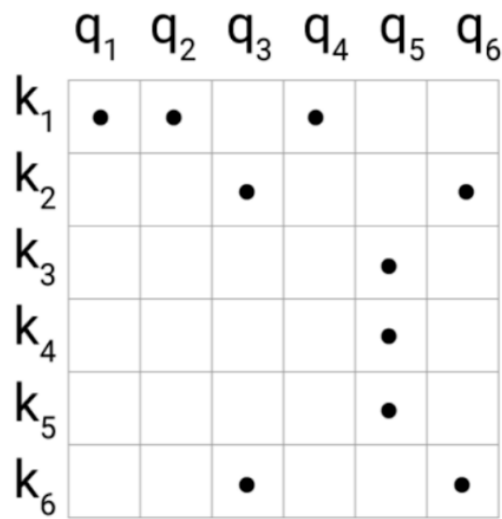


- Схема хэширования $x \mapsto h(x)$ является locality-sensitive, если сохраняет информацию о расстоянии между данными так, что близкие вектора получают одинаковые хэши, а далекие получают совершенно разные

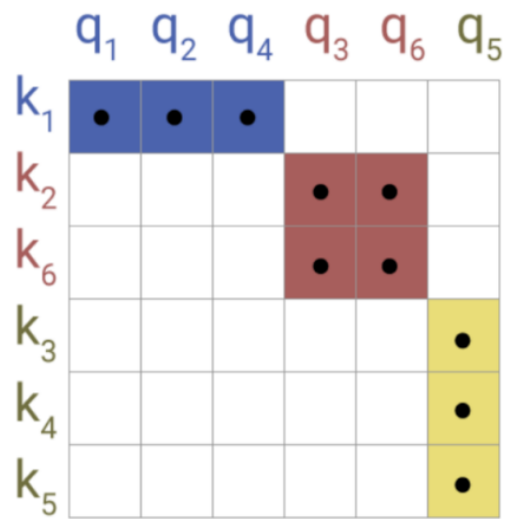
Locality-Sensitive Hashing Attention

- Реформер адаптирует схему хэширования
- Фиксированная случайная матрица $R \in \mathbb{R}^{d \times b/2}$, b - гиперпараметр
- $h(x) = \operatorname{argmax}([xR; -xR])$

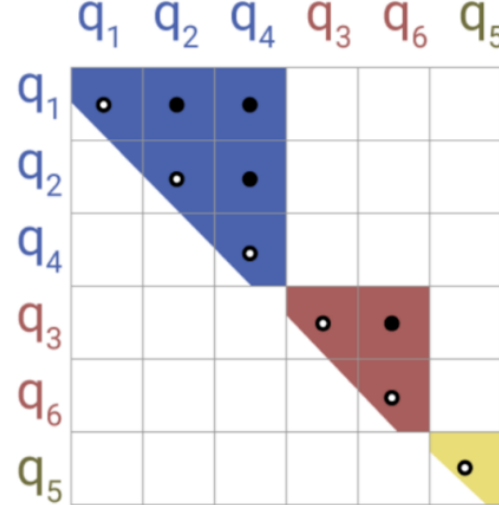
Locality-Sensitive Hashing Attention



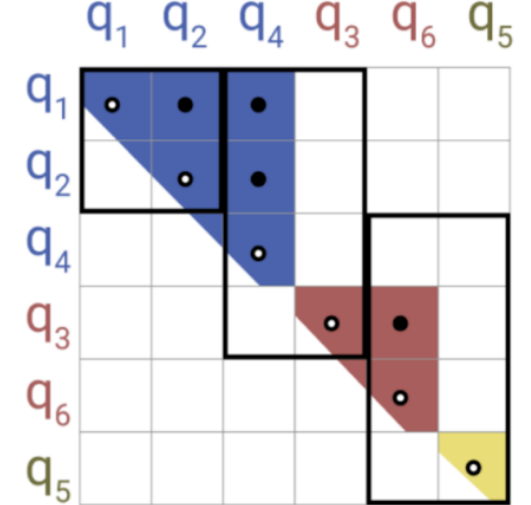
(a) Normal



(b) Bucketed



(c) $Q = K$

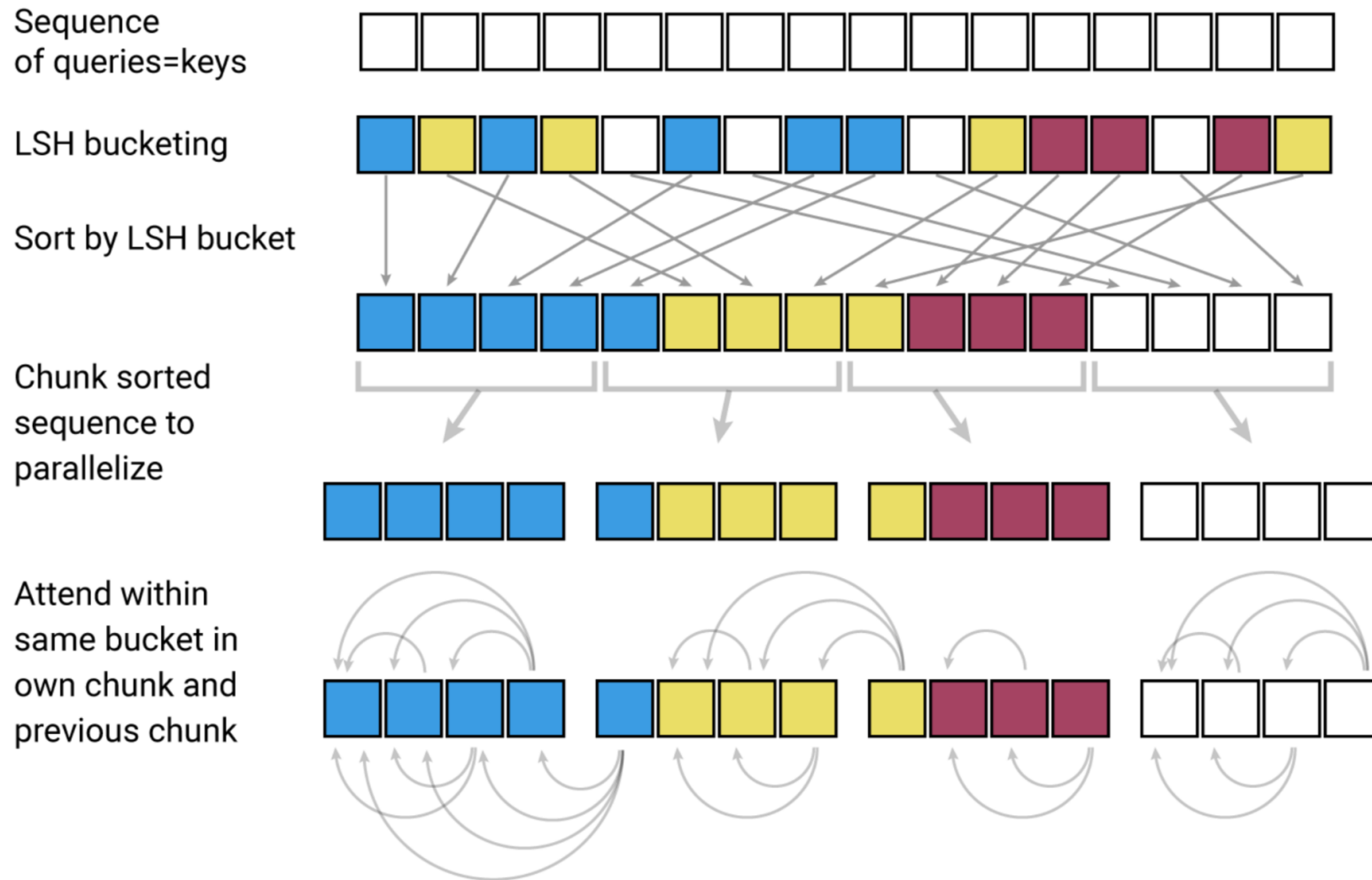


(d) Chunked

Locality-Sensitive Hashing Attention

- Запрос относится только к позициям из той же группы
- Матрица attention для full attention часто разреженная
- Сортируем ключи и запросы относительно их хэш групп.

Locality-Sensitive Hashing Attention



LSH attention:

- группировка,
- сортировка,
- разделение,
- **ВЫЧИСЛЕНИЕ** attention

Обратимые остаточные нейросети

- Слой $x \mapsto y$, normal residual layer: $y = x + F(x)$
- Обратимые слои делят input и output на пары $(x_1, x_2) \mapsto (y_1, y_2)$

Вычисляем:

$$y_1 = x_1 + F(x_2), y_2 = x_2 + G(y_1)$$

Обратить легко:

$$x_2 = y_2 - G(y_1), x_1 = y_1 - F(x_2)$$

Обратимые остаточные нейросети

- Реформер применяет ту же идею к трансформеру
- Комбинируем attention (F) и слои с прямой связью (G) внутри реверсивного блока:
- $Y1 = X1 + \text{Attention}(X2)$, $Y2 = X2 + \text{FeedForward}(Y1)$

Обратимые остаточные нейросети

- Память может быть еще уменьшена
- Разделим для этого вычисления с прямой связью:

$$Y_2 = [Y^{(1)}_2; \dots; Y^{(c)}_2] = \\ [X^{(1)}_2 + \text{FeedForward}(Y^{(1)}_1); \dots; X^{(c)}_2 + \text{FeedForward}(Y^{(c)}_1)]$$

- Полученный обратимый трансформер может не хранить активацию в каждом слое

Результаты:

- Благодаря LSH attention уменьшили сложность до $O(L \log L)$.
- Благодаря обратимости остаточных слоев храним активации только 1 раз за обучение, вместо N