

# Глава 1

## Алгоритмы визуализации данных на клиенте

В данной главе описываются реализации двух алгоритмов визуализации данных: Fruchterman Reingold и t-SNE.

Несмотря на наличие множества реализаций этих алгоритмов на JavaScript, ни одна из них не использует GPU. При этом нет простого способа портировать существующие реализации этих алгоритмов на GPU с других языков программирования, так как WebGL API очень ограничен по сравнению с OpenGL и другими API для графических карт.

### 1.1 WebGL для обработки данных

### 1.2 Fruchterman Reingold

Это классический алгоритм рисования графа из категории force-directed drawing.

Он использует физическую аналогию, в которой ребрам соответствуют пружины, а вершинам — одинаково заряженные частицы. Таким образом, соединенные ребром вершины стремятся быть на фиксированном расстоянии друг от друга, равно оптимальной длине пружины, а не соединенные вершины отталкиваются друг от друга.

На практике, закон Гука для силы пружины не используется, так как он слишком сильно действует на соединенные вершины, находящиеся в разных частях графа.

Для алгоритма Fruchterman Reingold силы притяжения и отталкивания равны  $f_a(d) = d^2/k$  и  $f_r(d) = -k^2/d$ , где  $k$  — оптимальная длина пружины, при которой эти силы сбалансированы.

Обычно, физическая симуляция осуществляется при помощи интегрирования Верле, но в этом алгоритме силы по сути являются скоростями, так как это ускоряет сходимость.

Алгоритм начинает со случайных позиций вершин и останавливается, когда система достигает равновесия. Для ускорения этого процесса вводится понятие *температуры*, которая ограничивает максимальное изменение

положения вершин и уменьшается по заданному закону.

### 1.3 t-SNE

Применение этого алгоритма считается одним из наиболее мощных методов уменьшения размерности до 1-3 измерений.

### 1.4 Схема вычислений.

Два приведенных алгоритма имеют общую структуру:

```
const edges: buffer[m] = serialize_edges()
positions: buffer[n, dim] = random_positions()
attractions, repulsions: buffer[n, dim]
grid: buffer[s]
finished: boolean = false
iteration: integer = 0

while not finished:
    iteration += 1
    attractions = calc_attractions(positions, edges)
    grid = calc_grid(positions)
    repulsions = calc_repulsions(positions, grid)
    positions, finished = calc_positions(
        positions,
        attractions,
        repulsions,
        iteration,
    )
```

Здесь  $n$  — число вершин,  $m$  — число ребер,  $dim$  — размерность пространства визуализации (от 1 до 3).

На вход алгоритму поступают два буфера: `edges` и `positions`. Позиции инициализируются случайными числами в заданных пределах. Ребра преобразуются в буфер при помощи одной из нескольких схем, перечисленных далее.

После этого следует некоторое число итераций. Для алгоритма Fruchterman Reingold номер итерации нужен для уменьшения температуры со временем, t-SNE останавливается, когда вершины перестают менять положение больше, чем на  $\epsilon$ .

Силы притяжения действуют только между парами вершин, соединенных ребром. Таким образом, вычисление сумм этих сил для всех вершин не возможно без прохода по всем ребрам, т.е. за  $O(m)$  времени и памяти. Улучшение скорости этой части возможно через уменьшение числа ребер.

Силы отталкивания действуют между всеми парами вершин, что позволяет делать некоторые оптимизации, при этом уменьшая точность этих вычислений на несколько процентов. Сложность этой части варьируется  $O(n \log n)$  до  $O(n^2)$  по времени и  $O(n)$  по памяти.

При использовании оптимизаций требуется дополнительная структура `grid`, которая, как правило, требует одного прохода по позициям вершин и  $O(n)$  памяти.

Последняя функция `calc_positions` также проходит по всем вершинам, при этом она может содержать дополнительную логику, которая проверяет, например, выход вершин за границы ограничивающего прямоугольника.

Исходя из приведенных выше оценок, лучше всего для исполнения на GPU подходят функции `calc_attractions` и `calc_repulsions`.

## 1.5 Выбор ребер и их представление

Вычисление сил притягивания определяется выбором подмножества ребер и способом их представления в виде текстуры.

Есть три способа выбора ребер:

- (a) Все ребра.
- (b) Оставить не более  $k$  ребер наибольшего веса у каждой вершины.
- (c) Оставить не более  $k$  ребер наибольшего веса у всего графа.

Также есть четыре способа представления ребер:

- (d) Матрица смежности. Элементы — веса ребер.
- (e) Матрица  $n \times 2k$ , в которой строки соответствуют вершинам, в каждой строке слева направо указан список ребер в виде чередующегося списка номеров вершин и весов ребер.  $k$  — наибольшая степень вершины.
- (f) Вектор длины  $n + 2m$ . Первые  $n$  чисел — «указатели» на индексы вектора, с которых начинается список ребер соответствующей вершины. Следующие  $2m$  чисел являются объединением списков ребер, выходящих из каждой вершины.
- (g) Матрица  $n' \times (1 + 2s)$ . Строки соответствуют вершинам, при этом одной вершине может соответствовать несколько подряд идущих строк. Номер вершины, соответствующей данной строке указан в первом столбце. В остальных столбцах указан список ребер и их весов. Если степень вершины больше  $s$ , она занимает несколько строк.

Способ (e) лучше всего подходит вместе со способом выбора ребер (b), иначе, если есть хотя бы одна вершина степени  $n - 1$ , то этот способ будет хуже по времени и памяти, чем (d).

Способ (f) требует меньше всего памяти, но, в связи с особенностью работы GPU, время обработки всех вершин будет пропорционально максимальному времени обработки одной вершины, то есть максимальной степени  $k$ . Если есть хотя бы одна вершина степени  $n - 1$ , то его время работы будет дольше, чем у способа (d).

Наконец, способ (g) требует не сильно больше памяти по сравнению с (f), при этом наиболее эффективно распараллеливается на GPU. Единственный недостаток состоит в том, что после его применения необходимо агрегировать результаты строк, соответствующих одной вершине.

## 1.6 Реализация вычисления сил притяжения

В данном разделе приводится псевдокод реализации функции `calc_attractions` для представления ребер ( $g$ ). Его несложно упростить для представления (e).

Для хранения информации о ребрах используется RGBA-текстуры вместо `floating point`. Это позволяет хранить два 16-битных числа в одном пикселе. Первое число используется для номера вершины в списке ребер (которых будет значительно меньше 65536), а второе — для веса ребра. Такое представление уменьшит точность хранения весов ребер, но, с учетом того, что веса ребер входят в формулы линейно и не связаны уравнениями с другими переменными, это не должно значительно ухудшить качество визуализации.

Значение 0 для веса ребер означает отсутствие ребра. Это может потребоваться, чтобы заполнить оставшиеся ячейки в строках представления ( $g$ ). Первая колонка матрицы, кодирующая номер вершины, соответствующей данной строке, использует только первое 16-битное число.

Результатом работы шейдера будут являться `dim` вещественных текстур  $1 \times n$ , содержащих вектора сил, действующие на вершины от ребер по строкам. На вход он получает текстуру `edges` размера  $n' \times (1 + s)$  и текстуру `positions` размера  $n \times dim$ .

Для каждой размерности `dim` код компилируется отдельно. Псевдокод шейдера для  $dim = 2$ :

```
const edges_chunks, s: integer;
const max_weight: float;

int high(rgba pixel):
    return pixel.r * 256 + pixel.g

int low(rgba pixel):
    return pixel.b * 256 + pixel.a

calc_attractions(chunk_no):
    int source = high(edges[chunk_no, 0])
    float x = positions[source, 0]
    float y = positions[source, 1]
    float x_res = 0
    float y_res = 0
    for i = 1 ... s:
        rgba edge_pixel = edges[chunk_no, i]
        int target = high(edge_pixel)
        float dx = positions[target, 0] - x
        float dy = positions[target, 1] - y
        float weight = low(edge_pixel) / 256 / 256 * max_weight
        float sqdist = dx ^ 2 + dy ^ 2
        float mul = calc_mul(sqdist, weight)
        x_res += mul * dx
        y_res += mul * dy
    return x_res, y_res
```

## **1.7 Подходы к вычислению сил отталкивания**

В обоих алгоритмах «узким местом» по времени выполнения является вычисление сил отталкивания между всеми парами вершин. Есть