

# Краткая история GPU и введение в OpenGL

Вычисления на видеокартах. Лекция 1

Полярный Николай

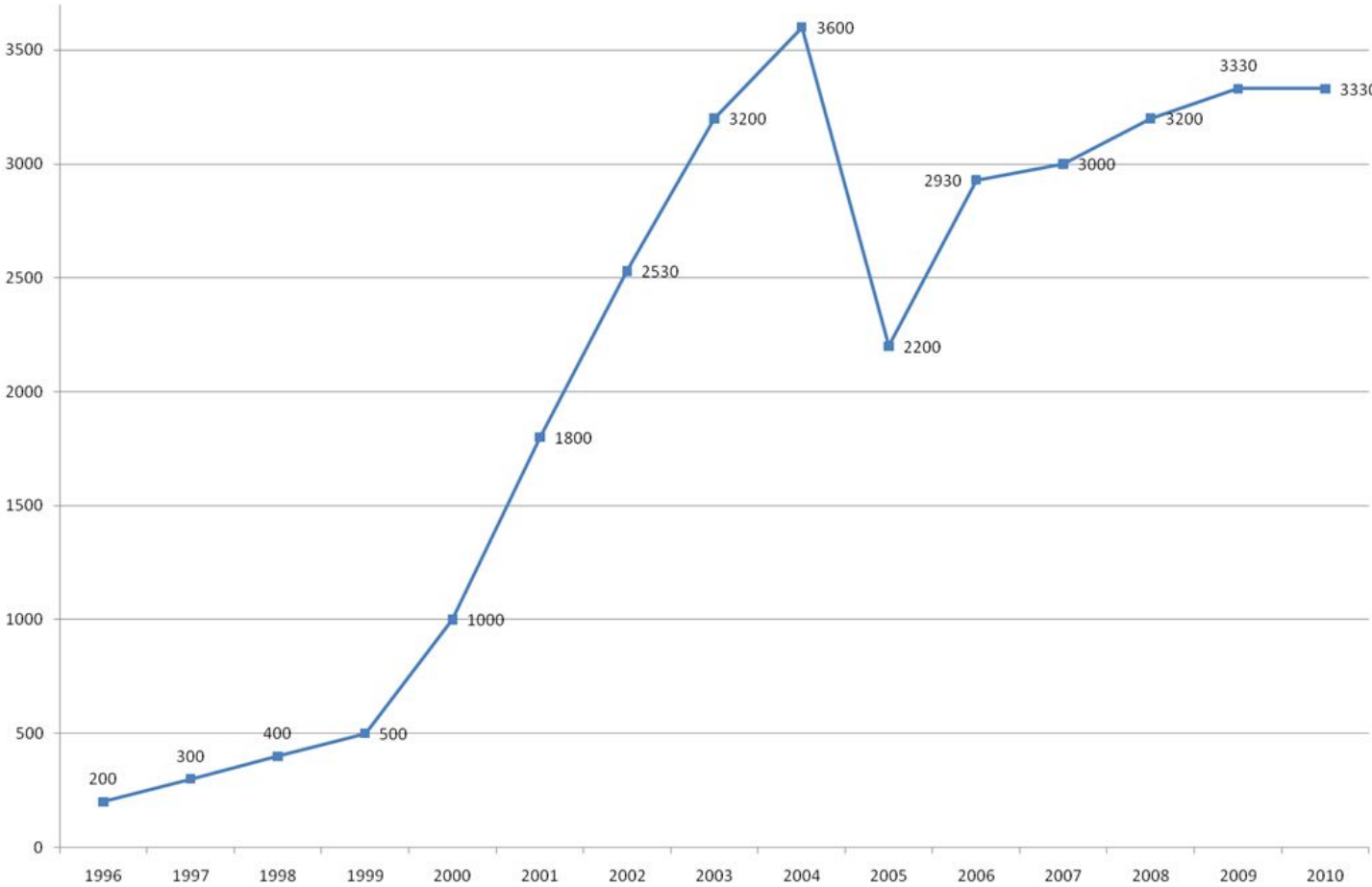
[polarnick239@gmail.com](mailto:polarnick239@gmail.com)

# Процессоры

Долго процессоры были одноядерными, и тактовая частота была их основной характеристикой.

В 2004 году рост частоты остановился на 3.6 ГГц (Pentium 4 Prescott).

Stock Clock Speed



# Процессоры

1. **Out-of-order** исполнение позволяет начинать заранее делать будущие команды (**Instruction-level parallelism**).

# Процессоры

1. **Out-of-order** исполнение позволяет начинать заранее делать будущие команды (**Instruction-level parallelism**).
2. Для улучшения ILP в случае ветвления кода существует **branch-prediction**.

# Процессоры

1. **Out-of-order** исполнение позволяет начинать заранее делать будущие команды (**Instruction-level parallelism**).
2. Для улучшения ILP в случае ветвления кода существует **branch-prediction**.
3. Часто узкое место не вычисления, а подгрузка данных из оперативной памяти, поэтому существуют **L1/L2/L3 кеш**и.

# Процессоры

1. **Out-of-order** исполнение позволяет начинать заранее делать будущие команды (**Instruction-level parallelism**).
2. Для улучшения ILP в случае ветвления кода существует **branch-prediction**.
3. Часто узкое место не вычисления, а подгрузка данных из оперативной памяти, поэтому существуют **L1/L2/L3 кеши**.
4. Скрывать задержки обращения к памяти помогает так же **Hyper-threading** и **Simultaneous multithreading**.

# Процессоры

1. **Out-of-order** исполнение позволяет начинать заранее делать будущие команды (**Instruction-level parallelism**).
2. Для улучшения ILP в случае ветвления кода существует **branch-prediction**.
3. Часто узкое место не вычисления, а подгрузка данных из оперативной памяти, поэтому существуют **L1/L2/L3 кеш**и.
4. Скрывать задержки обращения к памяти помогает так же **Hyper-threading** и **Simultaneous multithreading**.
5. **Многоядерность**. В не-серверном компьютере выросла до 4 ядер в 2008 году. В 2017 году - до 8-16 ядер.



# Процессоры: **Single instruction, multiple data**

- 1997, **MMX**: 64 бита, целочисленные операции, **8xchar, 4xshort, 2xint**
- 1999, **SSE**: 128 бит, floating-point операции, **4xfloat**
- 2000, **SSE2**: 128 бит, теперь и над целыми, **16xchar, 8xshort, 4xint, 2xlong, 4xfloat, 2xdouble**
- 2004-2006, **SSE3 / SSSE3 / SSE4 / SSE4.1 / SSE4.2** - расширения набора инструкций (горизонтальная работа с регистрами, специальные функции для видеокодирования и обработки строк)

# Процессоры: **Single instruction, multiple data**

- 2011, **AVX**: 256 бит, floating-point операции, **8xfloat, 4xdouble**
- 2013, **AVX2**: 256 бит, теперь и над целыми, **32xchar, 16xshort, 8xint, 4xlong, 8xfloat, 4xdouble**
- 2017, **AVX-512**: 512 бит, много подрасширений, мало где доступно, **16xfloat, 8xdouble**

# Множество Мандельброта без SIMD интринсик

```
float x = 0.0f;
float y = 0.0f;
for (iteration = 0; iteration < MAX_ITERATIONS; iteration++) {
    float xn = x * x - y * y + x0;
    y = 2 * x * y + y0;
    x = xn;
    if (x * x + y * y > INFINITY) {
        break;
    }
}
```

# Множество Мандельброта с SSE, SSE2, SSSE3

```
__m128i maskAll = _mm_setzero_si128();
__m128i iters = _mm_setzero_si128();
__m128 xs = _mm_set1_ps(0.0f);
__m128 ys = _mm_set1_ps(0.0f);
for (iteration = 0; iteration < MAX_ITERATIONS; iteration++) {
    __m128 xsn = _mm_add_ps(_mm_sub_ps(_mm_mul_ps(xs, xs), _mm_mul_ps(ys, ys)), xs);
    __m128 ysn = _mm_add_ps(_mm_mul_ps(_mm_mul_ps(_mm_set1_ps(2.0f), xs), ys), _mm_set1_ps(y0));
    xs = _mm_add_ps(_mm_andnot_ps((__m128) maskAll, xsn), _mm_and_ps((__m128) maskAll, xs));
    ys = _mm_add_ps(_mm_andnot_ps((__m128) maskAll, ysn), _mm_and_ps((__m128) maskAll, ys));

    maskAll = _mm_or_si128(_mm_castps_si128(_mm_cmpge_ps(_mm_add_ps(_mm_mul_ps(xs, xs), _mm_mul_ps(ys, ys),
    iters = _mm_add_epi16(iters, _mm_andnot_si128(maskAll, _mm_set1_epi16(1)));
    int mask = _mm_movemask_epi8(maskAll);
    if (mask == 0xffff) {
        break;
    }
}
iters = _mm_shuffle_epi8(iters, _mm_setr_epi8(12, 13, 8, 9, 4, 5, 0, 1, -1, -1, -1, -1, -1, -1, -1, -1));
```

# Множество Мандельброта с AVX и AVX2

```
__m256i maskAll = _mm256_setzero_si256();
__m256i iters = _mm256_setzero_si256();
__m256 xs = _mm256_set1_ps(0.0f);
__m256 ys = _mm256_set1_ps(0.0f);
for (iteration = 0; iteration < MAX_ITERATIONS; iteration++) {
    __m256 xsn = _mm256_add_ps(_mm256_sub_ps(_mm256_mul_ps(xs, xs), _mm256_mul_ps(ys, ys)), xs0);
    __m256 ysn = _mm256_add_ps(_mm256_mul_ps(_mm256_mul_ps(_mm256_set1_ps(2.0f), xs), ys), _mm256_set1_ps
xs = _mm256_add_ps(_mm256_andnot_ps((__m256) maskAll, xsn), _mm256_and_ps((__m256) maskAll, xs));
ys = _mm256_add_ps(_mm256_andnot_ps((__m256) maskAll, ysn), _mm256_and_ps((__m256) maskAll, ys));

    maskAll = (__m256i) _mm256_or_ps(_mm256_cmp_ps(_mm256_add_ps(_mm256_mul_ps(xs, xs), _mm256_mul_ps(ys,
iters = _mm256_add_epi16(iters, _mm256_andnot_si256(maskAll, _mm256_set1_epi16(1)));
    int mask = _mm256_movemask_epi8(maskAll);
    if (mask == (int) 0xffffffff) {
        break;
    }
}
iters = _mm256_shuffle_epi8(iters, _mm256_setr_epi8(0, 1, -1, -1, 4, 5, -1, -1, 8, 9, -1, -1, 12, 13, -1,
```

# Software 3D

Вплоть **до 1996** года даже игры вроде DOOM, Quake и Duke Nukem 3D обсчитывали визуализацию 3D пространства на процессоре. Процессор справлялся в первую очередь благодаря низкому расширению экрана 320x200 и относительно простой графике.



# Специализированные 3D-ускорители

Около **1995-1997** начали набирать популярность специализированные 3D-ускорители (такие как S3 ViRGE, ATI 3D Rage, 3dfx Voodoo) способные отрисовывать более сложную графику гораздо эффективнее.



# 3D-ускорители: программируемые шейдеры

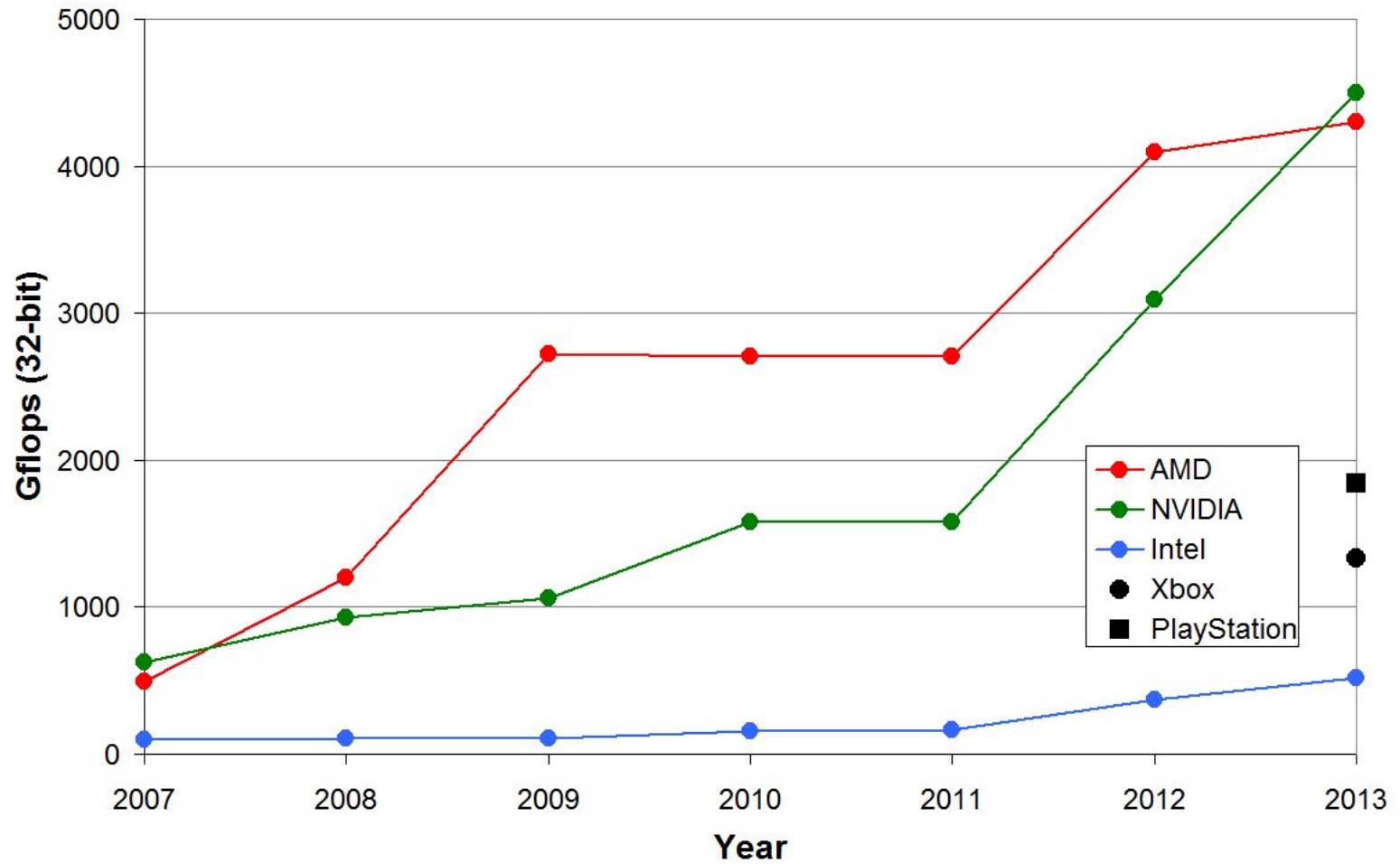
Спецэффекты в играх становились все сложнее и стало невозможно продолжать предлагать программистам лишь жестко ограниченный набор 3D функциональности, возникла потребность дать гибкий инструмент для программирования видеокарты. Так появились программируемые шейдеры.

Стало возможно в каждом пикселе посчитать произвольную математику, и это позволило реализовать много интересных спецэффектов (таких как реалистичная вода, сложные модели освещения и т.п.).

# Зарождение GPGPU

Благодаря популярности игр, архитектурным особенностям видеокарт (задача отрисовки обладает массовым параллелизмом, т.к. пиксели на экране можно обсчитывать независимо) и все возрастающей сложности игровой графики производительность видеокарт росла гораздо быстрее процессоров:

## GPU and CPU Peak Performance Trends



# Зарождение GPGPU

Поэтому видеокарты специализирующиеся на отображении 3D-графики начали использовать и для других задач, таких как обсчет физики и моделирования сложных процессов.

Это стало возможно благодаря гибкости программируемых шейдеров, но тем не менее из-за специализации на 3D-графике гибкость шейдеров была ограничена, и для general-purpose задач возникла потребность нового, еще более гибкого API.

# Появление GPGPU API

Стали появляться API лучше подходящие для general-purpose вычислений:

- 2006, **Close to Metal** - поддерживался только на видеокартах от **AMD**, позже был заменен на **OpenCL**
- 2007, **CUDA** - поддерживается только на видеокартах от **NVidia**
- 2009, **OpenCL** - открытый стандарт поддерживающийся практически везде

# Появление GPGPU API

Так же позже появилось несколько узко специализированных API:

- 2011, RenderScript - специально для Android, изначально ставил целью заменить OpenGL, с 2013 года стал пытаться заменить OpenCL

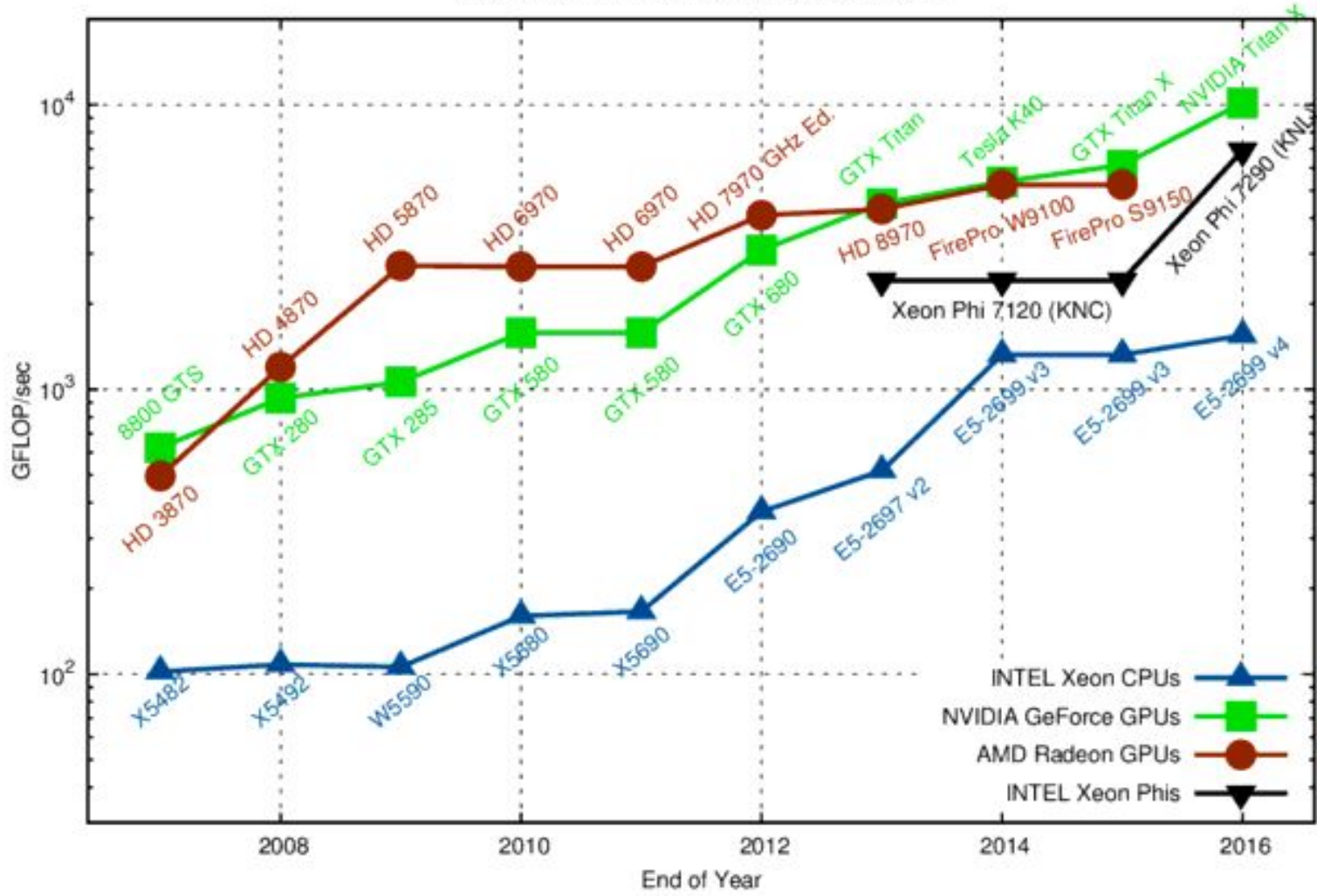
- 2014, Metal - специально для устройств Apple, и графика, и вычисления

В результате на сегодняшний день для вычислений общего назначения обладающих массовым параллелизмом на видеокартах существуют удобные API, и реализовать потенциал видеокарты с их помощью гораздо проще чем теоретическую мощность процессора с помощью SIMD.

# Производительность GPU

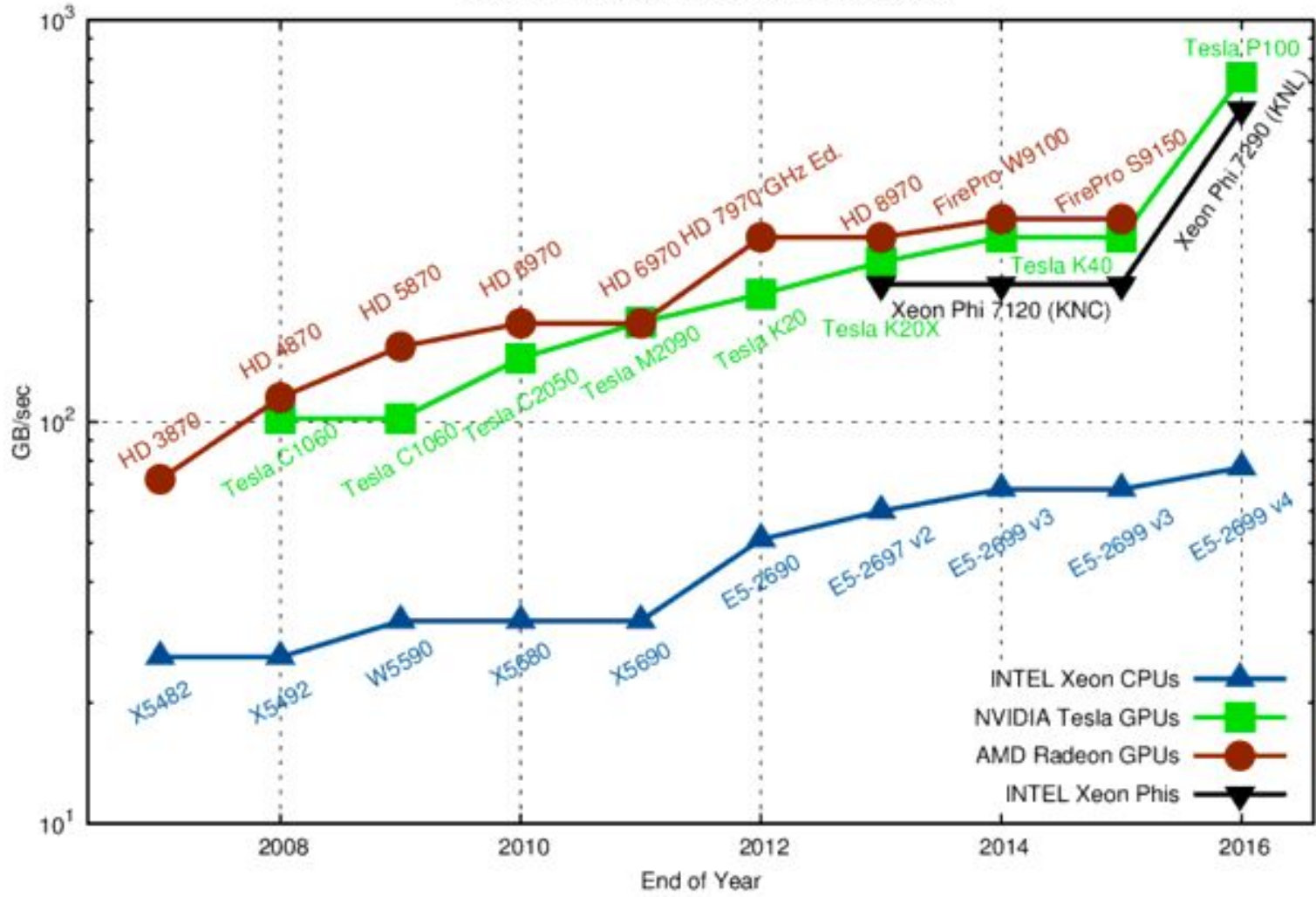
Производительность видеокарт же на порядок выше по обоим критическим показателям - вычислительная мощность и пропускная способность видеопамати:

Theoretical Peak Performance, Single Precision





Theoretical Peak Memory Bandwidth Comparison



# Мотивация

Цель курса:

- Понять что такое **массовый параллелизм** и как спроектировать алгоритм с его учетом
- Научиться использовать **OpenCL**
- Понять почему видеокарты хорошо подходят только для задач обладающих **массовым параллелизмом**
- Понять общие свойства архитектуры видеокарт и тем самым научиться оптимизировать алгоритмы
- Когда-нибудь потом, делая исследовательский проект, внезапно понять что вы можете за один вечер ускорить вычисления в десятки раз, что существенно увеличит скорость экспериментирования

# OpenCL платформы и устройства

## Platforms:

INTEL  
platform

AMD  
platform

NVIDIA  
platform

...

## Devices:

Intel HD  
(Integrated GPU)

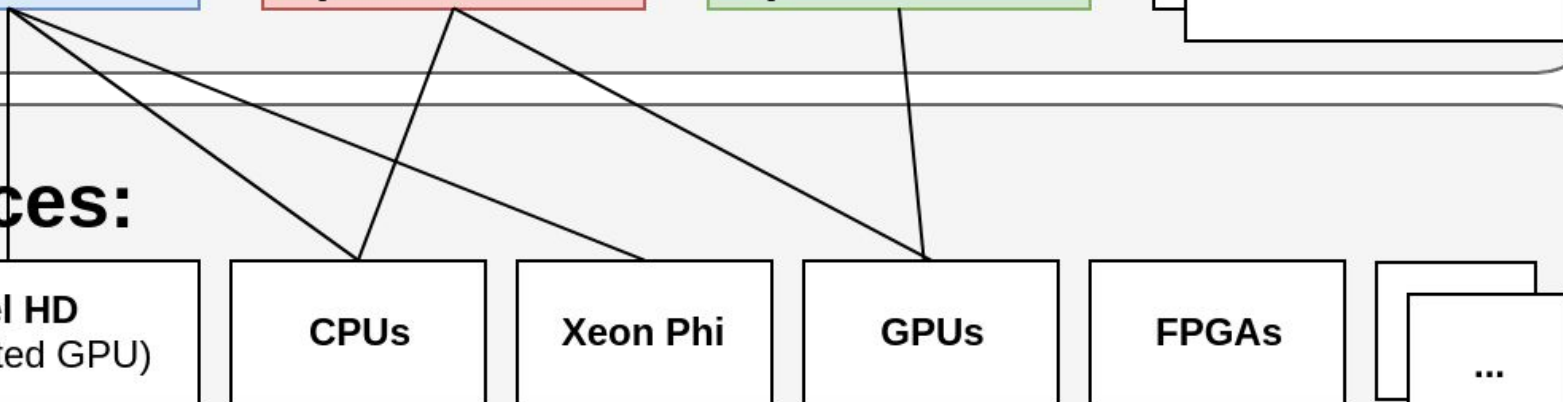
CPUs

Xeon Phi

GPUs

FPGAs

...



# Модель вычислений

