

Job scheduler

просто, насколько можно

но не проще

Арсений “Zeux” Капулкин

CREAT Studios

a.kapoulkine@creatstudios.com

Free lunch is over!

- The biggest sea change in software development since the OO revolution is knocking at the door, and its name is Concurrency.

– *Herb Sutter*

Параллельный мир

- Много ядер общего назначения
 - PS3: 6 SPU
 - Xbox360: 3 PowerPC ядра
- Разные подходы
 - Параллельность по потокам (thread-parallel)
 - Параллельность по задачам (task-parallel)
 - Параллельность по данным (data-parallel)

Термины

- Job
 - Код, производящий преобразование над данными
- Batch
 - Ссылка на job + данные
- Scheduler
 - Система, выполняющая batch-и

Job

- Функция, обрабатывающая данные
 - Влияние ограничено входными данными
 - Нет доступа к глобальным переменным
- Выполняется много раз на кадр
 - Возможно, одновременно
 - Возможно, на одних и тех же данных
- Нет preemption
 - Требуется дополнительной памяти

Scheduler v1.0

- Максимально простой код
 - Достаточно удобный
 - Быстрый
- Добавление батчей на выполнение
- Ожидание выполнения
 - Синхронизация
 - Два уровня: батч и группа батчей
 - 32 группы батчей, uint32 счетчик на группу

Scheduler v1.0

- Глобальная очередь батчей
 - Фиксированный размер очереди
 - Lock-free
 - Операции push и pop
 - Операции копируют данные
 - Устраняет ручной менеджмент батчей
 - Операций empty, size нет
 - Не имеют смысла в многопоточной среде
 - Применима не только в schedulers!

Lock-free очередь

- Очередь (FIFO) – важный примитив
- Много опубликованных реализаций
 - MPMC – Multi-Producer, Multi-Consumer
 - Implementing Lock-Free Queues [94] – RACE!
 - Correction of a Memory Management Method for Lock-Free Data Structures [95]
 - Optimised Lock-Free FIFO Queue [01] – RACE!
 - Optimised Lock-Free FIFO Queue [03]
 - Optimized Lock-Free FIFO Queue continued [05]

Lock-free очередь – 1/3

```
1. template <class T>
2. class TLockFreeQueue
3. {
4.     struct TListNode
5.     {
6.         TListNode *volatile Next;
7.         T Data;
8.     };
9.
10.    struct TRootNode
11.    {
12.        TListNode *volatile PushQueue;
13.        TListNode *volatile PopQueue;
14.        TListNode *volatile ToDelete;
15.        TRootNode *volatile NextFree;
16.
17.        TRootNode() : PushQueue(0), PopQueue(0),
18.            ToDelete(0), NextFree(0) {}
19.    };
20.
21.    static void EraseList(TListNode *n)
22.    {
23.        while (n) {
24.            TListNode *keepNext = n->Next;
25.            delete n;
26.            n = keepNext;
27.        }
28.    }
29.
30.    TRootNode *volatile JobQueue;
31.    volatile long FreeMemCounter;
32.    TRootNode *volatile FreePtr;
33.
34.    void TryToFreeAsyncMemory()
35.    {
36.        TRootNode *current = FreePtr;
37.        if (current == 0)
38.            return;
39.        if (atomic_add(&FreeMemCounter, 0) == 1) {
40.            // we are the last thread, try to cleanup
41.            if (cas(&FreePtr, (TRootNode*)0, current)) {
42.                // free list
43.                while (current) {
44.                    TRootNode *p = current->NextFree;
45.                    EraseList(current->ToDelete);
46.                    delete current;
47.                    current = p;
48.                }
49.            }
50.        }
51.    }
52.
53.    void AsyncRef()
54.    {
55.        atomic_add(&FreeMemCounter, 1);
56.    }
57.
58.    void AsyncUnref()
59.    {
60.        TryToFreeAsyncMemory();
61.        atomic_add(&FreeMemCounter, -1);
62.    }
63.
64.    void AsyncDel(TRootNode *toDelete, TListNode *lst)
65.    {
66.        toDelete->ToDelete = lst;
67.        // we are the last thread, try to cleanup
68.        if (cas(&FreePtr, toDelete, toDelete->NextFree))
69.            break;
70.    }
71.
72.    void AsyncUnref(TRootNode *toDelete, TListNode *lst)
73.    {
74.        TryToFreeAsyncMemory();
75.        if (atomic_add(&FreeMemCounter, -1) == 0) {
76.            // no other operations in progress, can safely
77.            // reclaim memory
78.            EraseList(lst);
79.            delete toDelete;
80.        } else {
81.            // Dequeue() is in progress, put node to free
82.            list
83.            AsyncDel(toDelete, lst);
84.        }
85.    }
86.
87.    struct TListInverter
88.    {
89.        TListNode *Copy;
90.        TListNode *Tail;
91.        TListNode *PrevFirst;
92.
93.        TListInverter() : Copy(0), Tail(0), PrevFirst(0) {}
94.        ~TListInverter() {}
95.
96.        if (Copy)
97.            Copy = Copy;
98.        EraseList(Copy);
99.    }
100.
101.    void CopyWasUsed()
102.    {
103.        Copy = 0;
104.        Tail = 0;
105.        PrevFirst = 0;
106.    }
107.
108.    void DoCopy(TListNode *ptr)
```

Lock-free очередь – 2/3

```
103.         {
104.             TListNode *newFirst = ptr;
105.             TListNode *newCopy = 0;
106.             TListNode *newTail = 0;
107.             while (ptr) {
108.                 if (ptr == PrevFirst) {
109.                     // short cut, we have copied
110.                     this part already
111.                     Tail->Next = newCopy;
112.                     newCopy = Copy;
113.                     Copy = 0; // do not destroy prev
114.                     try
115.                     if (!newTail)
116.                         newTail = Tail; // tried to
117.                         invert same list
118.                         break;
119.                     }
120.                     TListNode *newElem = new TListNode;
121.                     newElem->Data = ptr->Data;
122.                     newElem->Next = newCopy;
123.                     newCopy = newElem;
124.                     ptr = ptr->Next;
125.                     if (!newTail)
126.                         newTail = newElem;
127.                     }
128.                     EraseList(Copy); // copy was useless
129.                     Copy = newCopy;
130.                     PrevFirst = newFirst;
131.                     Tail = newTail;
132.                 }
133.             };
134.             TLockFreeQueue(const TLockFreeQueue&) {}
135.             void operator=(const TLockFreeQueue&) {}
136.             public:
137.             TLockFreeQueue() : JobQueue(new TRootNode),
138.             FreemanCounter(0), FreePtr(0) {}
139.             ~TLockFreeQueue()
140.             {
141.                 AsyncRef();
142.                 AsyncUnref();
143.                 EraseList(JobQueue->PushQueue);
144.                 EraseList(JobQueue->PopQueue);
145.                 delete JobQueue;
146.             }
147.             void Enqueue(const T &data)
148.             {
149.                 TListNode *newNode = new TListNode;
150.                 newNode->Data = data;
151.                 TRootNode *newRoot = new TRootNode;
152.                 AsyncRef();
153.                 newRoot->PushQueue = newNode;
154.                 TRootNode *curRoot = JobQueue;
155.                 newRoot->PushQueue = newNode;
156.                 TListNode *tail = curRoot->PopQueue;
157.                 if (tail) {
158.                     // has elems to pop
159.                     if (!newRoot)
160.                         newRoot = new TRootNode;
161.                     newRoot->PushQueue = curRoot->
162.                     >PushQueue;
163.                     newRoot->PopQueue = tail->Next;
164.                     ASSERT(curRoot->PopQueue == tail);
165.                     if (curRoot->PopQueue == tail) {
166.                         *data = tail->Data;
167.                         tail->Next = 0;
168.                         AsyncUnref(curRoot, tail);
169.                         return true;
170.                     }
171.                     continue;
172.                 }
173.                 if (curRoot->PushQueue == 0) {
174.                     delete newRoot;
175.                     AsyncUnref();
176.                     return false; // no elems to pop
177.                 }
178.                 if (!newRoot)
179.                     newRoot = new TRootNode;
180.                 newRoot->PushQueue = 0;
181.                 listInverter.DoCopy(curRoot->PushQueue);
182.                 newRoot->PopQueue = listInverter.Copy;
183.                 ASSERT(curRoot->PopQueue == 0);
184.                 if (curRoot->PushQueue == newRoot, curRoot) {
185.                     newRoot = 0;
186.                     listInverter.CopyWasUsed();
187.                     AsyncDel(curRoot, curRoot->
188.                     >PushQueue);
189.                 }
190.             }
191.         }
192.     }
193.     bool Dequeue(T *data)
194.     {
195.         TRootNode *newRoot = 0;
196.         TListInverter listInverter;
197.         AsyncRef();
198.         for(;;) {
199.             TRootNode *curRoot = JobQueue;
200.             TListNode *tail = curRoot->PopQueue;
201.             if (tail) {
202.                 // has elems to pop
203.                 if (!newRoot)
204.                     newRoot = new TRootNode;
205.                 newRoot->PushQueue = curRoot->
206.                 >PushQueue;
207.                 newRoot->PopQueue = tail->Next;
208.                 ASSERT(curRoot->PopQueue == tail);
209.                 if (curRoot->PopQueue == tail) {
210.                     *data = tail->Data;
211.                     tail->Next = 0;
212.                     AsyncUnref(curRoot, tail);
213.                     return true;
214.                 }
215.                 continue;
216.             }
217.             if (curRoot->PushQueue == 0) {
218.                 delete newRoot;
219.                 AsyncUnref();
220.                 return false; // no elems to pop
221.             }
222.             if (!newRoot)
223.                 newRoot = new TRootNode;
224.             newRoot->PushQueue = 0;
225.             listInverter.DoCopy(curRoot->PushQueue);
226.             newRoot->PopQueue = listInverter.Copy;
227.             ASSERT(curRoot->PopQueue == 0);
228.             if (curRoot->PushQueue == newRoot, curRoot) {
229.                 newRoot = 0;
230.                 listInverter.CopyWasUsed();
231.                 AsyncDel(curRoot, curRoot->
232.                 >PushQueue);
233.             }
234.         }
235.     }
236. }
```

Lock-free очередь – 3/3

```
202.         } else {
203.             newRoot->PopQueue = 0;
204.         }
205.     }
206. }
207. bool IsEmpty()
208. {
209.     AsyncRef();
210.     TRootNode *curRoot = JobQueue;
211.     bool res = curRoot->PushQueue == 0 &&
curRoot->PopQueue == 0;
212.     AsyncUnref();
213.     return res;
214. }
```

Lock-free очередь – 3/3

```
202.         } else {  
203.             newRoot->PopQueue = 0;  
204.         }  
205.     }  
206. }  
207. bool IsEmpty()  
208. {  
209.     AsyncRef();  
210.     TRootNode *curRoot = JobQueue;  
211.     bool res = curRoot->PushQueue == 0 &&  
212.         curRoot->PopQueue == 0;  
213.     AsyncUnref();  
214.     return res;  
215. }
```

- Бывают ли lock-free алгоритмы без багов?
 - Бывают.
 - А если кода – 4 экрана?
 - В этом коде была найдена утечка памяти.

Batch

```
struct Batch
```

```
{
```

```
    /* заголовок: 16 байт */
```

```
    ...
```

```
    /* данные для job-а: 112 байт */
```

```
    char user_data[112];
```

```
};
```

Queue

```
struct Queue
{
    Batch batches[QUEUE_SIZE];

    volatile uint32_t get;
    volatile uint32_t put;
};
```

Queue: принцип работы

- Кольцевой буфер
 - get – первый еще не обработанный батч
 - put – первый еще не записанный батч
 - $get \neq put$ – в буфере есть батчи
- get/put переполняются
 - Индексы батчей по модулю **QUEUE_SIZE**
 - **QUEUE_SIZE** – степень двойки
 - $2^{32} * 3000$ тактов = 1.2 часа

Queue: принцип работы

- У каждого батча есть два индекса
 - Глобальный – `uint32_t`, уникальный
 - Локальный – от 0 до `QUEUE_SIZE` - 1
- Операции начинаются с выбора батча
 - 1 атомарная инструкция
- Race conditions – в пределах 1 батча
 - ... или между батчами с совпадающим локальным индексом

Queue: добавление батча

```
// атомарный аналог index = q.put++  
uint32_t index = atomic_increment(&q.put);  
  
batch = &q.batches[index % QUEUE_SIZE];  
batch->user_data = user_data;
```

Queue: удаление батча

```
uint32_t index;
```

```
do {
```

```
    index = q.get;
```

```
    if (index == q.put) return QUEUE_EMPTY;
```

```
} while (!atomic_cas(&q.get, index, index + 1));
```

```
*result = q.batches[index % QUEUE_SIZE];
```

Race condition №1

```
// атомарный аналог index = q.put++  
uint32_t index = atomic_increment(&q.put);  
  
// index указывает на еще не заполненный  
// элемент; прочитанный batch может быть  
// целиком или частично устаревшим  
  
batch = &q.batches[index % QUEUE_SIZE];  
batch->user_data = user_data;
```

Race condition №1 – решение

```
// queue_push
```

```
batch->user_data = user_data;
```

```
memory_barrier();
```

```
batch->ready = true;
```

```
// queue_pop
```

```
while (batch->ready == false) yield();
```

```
*result = *batch;
```

Race condition №2

```
// атомарный аналог index = q.put++  
uint32_t index = atomic_increment(&q.put);
```

```
// index указывает на уже заполненный, но  
// еще не обработанный элемент; batch  
// будет целиком или частично утерян
```

```
batch = &q.batches[index % QUEUE_SIZE];  
batch->user_data = user_data;
```

Race condition №2 – решение?

```
// queue_push
```

```
while (batch->ready == true) yield();
```

```
// queue_pop
```

```
while (batch->ready == false) yield();
```

```
*result = *batch;
```

```
batch->ready = false;
```

Race condition №3

```
// queue_push
```

```
while (batch->ready == true) yield();
```

- Много потоков ждут одного батча
 - Один и тот же локальный индекс
 - Случается при переполнении очереди
- Race при освобождении батча
 - Много потоков модифицируют один батч

Race condition №3

- Требуется различать потоки, ждущие освобождения батча
 - Нужен уникальный идентификатор...
 - У нас есть глобальный индекс батча!
- Записать батч может поток с наименьшим глобальным индексом
 - Остальные ждут, пока этот батч не освободится

Batch

```
struct Batch
{
    /* заголовок: 12 байт + 4 байта padding */
    uint32_t index;
    struct Job job;

    /* данные для job-а: 112 байт */
    char user_data[112];
};
```

Queue: добавление батча

```
uint32_t index = atomic_increment(&q.put);  
batch = &q.batches[index % QUEUE_SIZE];  
while (batch_busy(index, batch->index)) yield();
```

```
batch->job = job;  
batch->user_data = user_data;  
memory_barrier();  
batch->index = index;
```

Queue: удаление батча

```
uint32_t index;  
do {  
    index = q.get;  
    if (index == q.put) return QUEUE_EMPTY;  
    batch = &q.batches[index % QUEUE_SIZE];  
    if (batch->index != index) { yield(); continue; }  
} while (!atomic_cas(&q.get, index, index + 1));  
*result = *batch;  
batch->index = index + 1; // освобождение
```

Queue: статус батча

- `batch->index == index`
 - Батч с глобальным номером `index` добавлен в очередь, но еще не удален
- `(index - batch->index) >= QUEUE_SIZE`
 - Батч с глобальным номером `!= index` добавлен в очередь – `batch_busy`
- В противном случае
 - Батч с глобальным номером `index` удален

Queue: статус батча

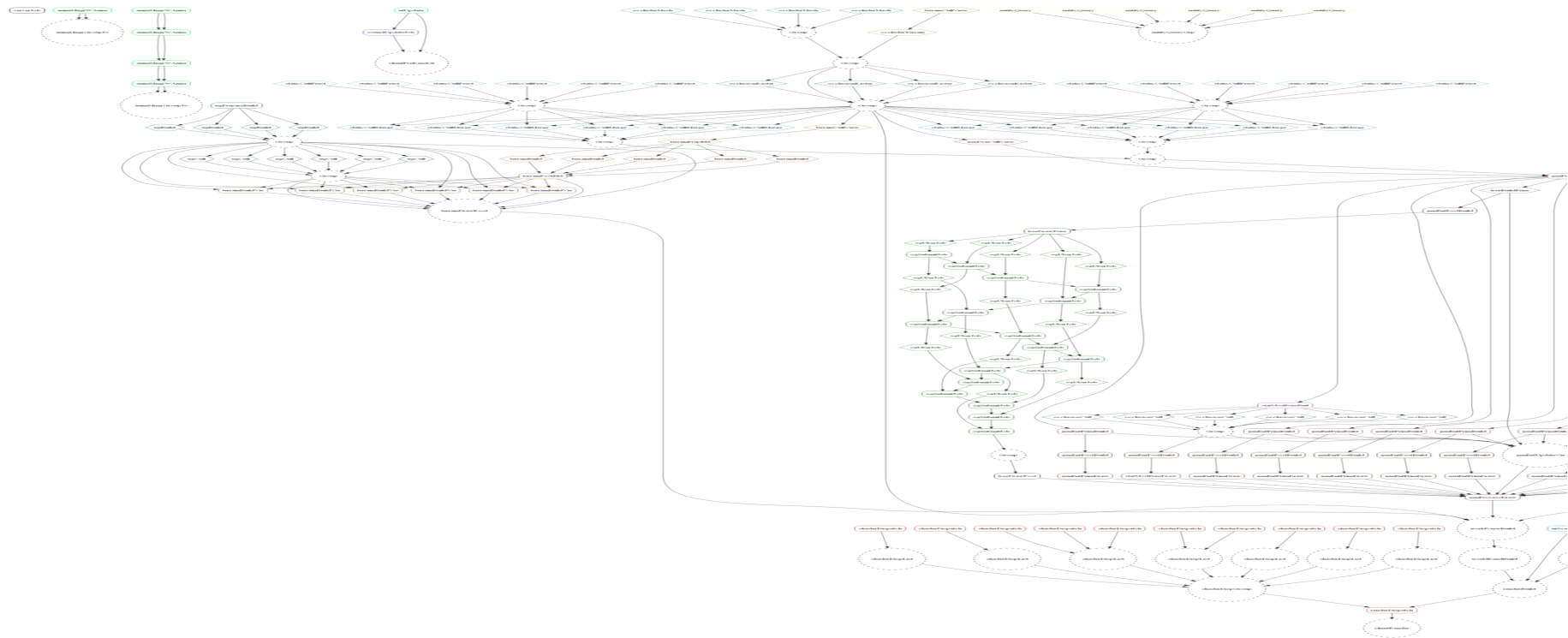
- `batch->index == index`
 - Батч с глобальным номером `index` добавлен в очередь, но еще не удален
- Отложим **освобождение** батча до конца его выполнения
 - Бесплатный **`wait_for_batch!`**

```
batch = &q.batches[index % QUEUE_SIZE];  
while (batch->index == index) yield();
```

Scheduler v1.0 - результаты

- Интерфейс для **task-parallel** задач
 - **uint32_t** push_batch
 - wait_for_batch
 - wait_for_group
- 1 атомарная операция на push/pop
 - +1 атомарная операция на group counter
- Простая реализация

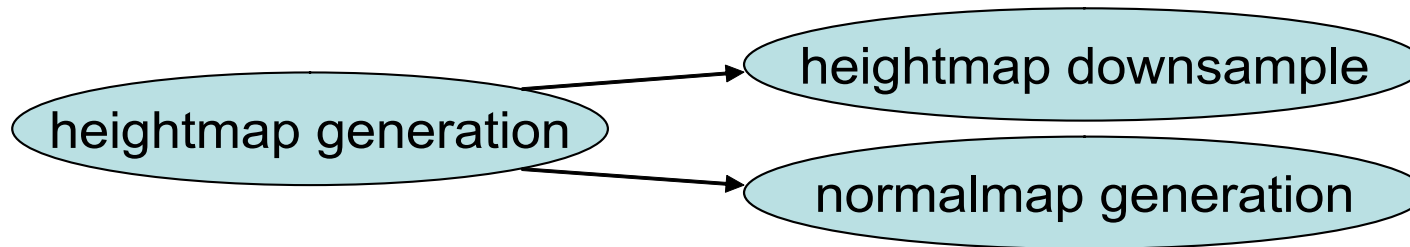
Граф зависимости батчей



(EA DICE Frostbite)

Граф зависимости батчей

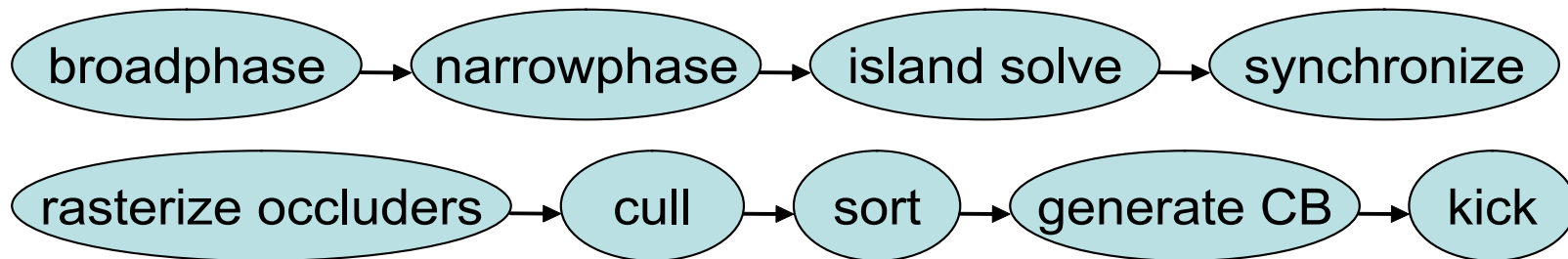
- Динамическое построение графа
 - Job запускает несколько новых батчей



- Эффективнее, чем добавлять все батчи и расставлять зависимости
 - Working set scheduler-а минимален

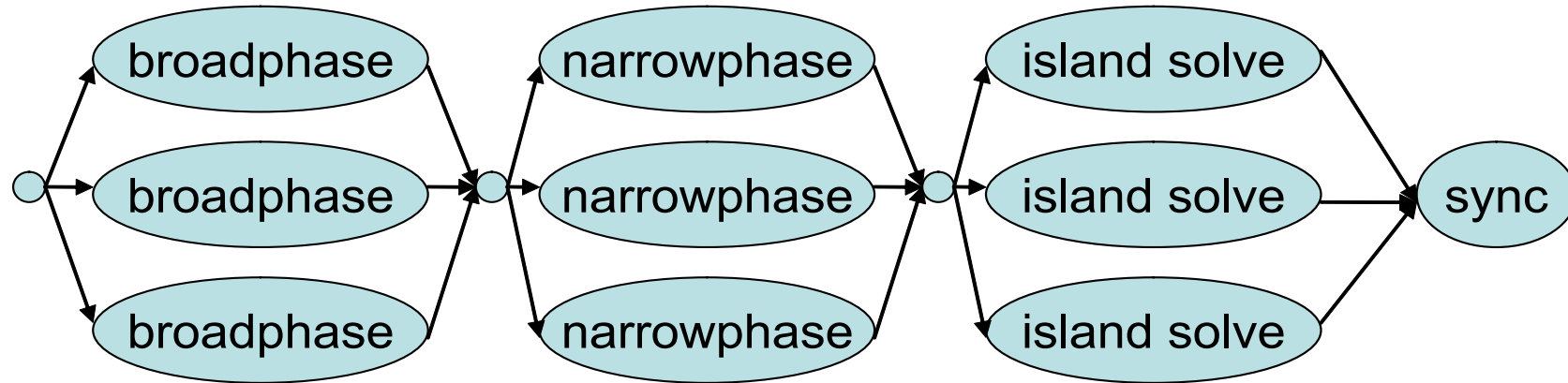
Граф зависимости батчей

- Параллелизм по данным
 - Частый паттерн №1: набор стадий
 - Возможно несколько одновременно работающих цепочек
 - Несколько проходов рендера



Граф зависимости батчей

- Параллелизм по данным
 - Частый паттерн №2: data-parallel стадии
 - Гораздо лучшая масштабируемость



Scheduler v2.0

- Data-parallel jobs
 - Система должна остаться простой
 - Ищем вдохновение в GPGPU
- CUDA / DirectCompute / OpenCL
 - Запуск batch-а с одним набором параметров на N ядрах
 - Каждому ядру приходят данные и индекс в блоке
 - blockIdx, blockDim

Batch block

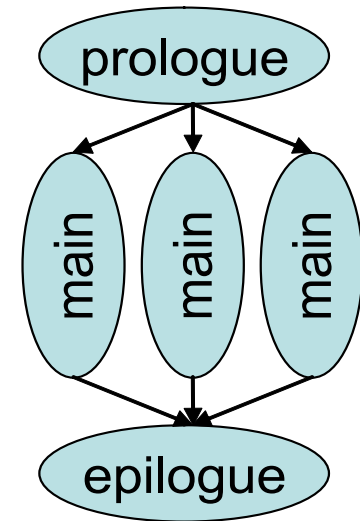
- В batch добавляются параметры блока
 - index – индекс батча в блоке [0..count)
 - count – общее число батчей в блоке
- Пример: parallel for
 - start = size / block.count * block.index
 - end = size / block.count * (block.index + 1)
 - for (i = start; i < end; ++i)
 - // обработка элемента номер i

Batch block - синхронизация

- Как синхронизировать job-ы в блоке?
 - GPU: `__syncthreads`
 - Требуется preemption
 - Требуется жестких критериев schedul-инга
 - Идея!
 - Пролог и эпилог

Пролог и эпилог

- Пролог
 - Выполняется один раз, до первого main
- Эпилог
 - Выполняется один раз, после последнего main
- Набор данных один и тот же
 - Для всех трех функций



Batch block - примеры

- Динамический parallel for
 - element = **atomic_increment**(&g_counter)
 - Block index используется для менеджмента scratch памяти
- Data-parallel for
 - Пролог: подготовка данных (если нужно)
 - main: обработка данных
 - Эпилог: запуск batch-а следующей стадии

Batch

```
struct Batch {  
    /* заголовок: 16 байт */  
    uint32_t index;  
    struct Block block;  
    struct Job job;  
  
    /* данные для job-а: 112 байт */  
    char user_data[112];  
};
```


Block

```
struct Block  
{  
    volatile uint16_t semaphore;  
    uint16_t count;  
};
```

Queue

```
struct Queue
{
    Batch batches[QUEUE_SIZE];

    volatile uint32_t get;
    volatile uint32_t get_block;
    volatile uint32_t put;
};
```

Queue: удаление батча

```
uint32_t index, block, new_index, new_block;  
do {  
    index, block = q.get, q.get_block;  
    if (index == q.put) return QUEUE_EMPTY;  
    batch = &q.batches[index % QUEUE_SIZE];  
    if (batch->index != index) { yield(); continue; }  
    bool last_batch = block + 1 == batch->count;  
    new_index = last_batch ? index + 1 : index;  
    new_block = last_batch ? 0 : block + 1;  
} while (!atomic_cas2(&q.get, index, new_index,  
    &q.get_block, block, new_block));
```

Queue: статус батча

- Батч остается занят, пока все job-ы не отработали
 - Включая пролог и эпилог
- Можно использовать тело батча для синхронизации между job-ами!
 - `volatile uint16_t` semaphore;
 - Batch занимает cache line на PS3/XBox360
 - 128b размер и выравнивание

Обработка пролога

- Если пролог есть:
 - Если `batch.index == 0`, увеличить семафор
 - В противном случае ждать, пока семафор не станет 1
- Наличие пролога означает ожидание
 - Применяется редко
 - Обычно лучше вынести пролог в отдельный батч

Обработка эпилога

- Увеличить семафор на 1
- Если семафор был равен `block.count`, то это последний батч
 - В случае, если пролога нет, семафор будет равен `block.count - 1`
 - Запустить эпилог, если нужно
 - Освободить батч

Дополнительный overhead

- Батч с `block.count == 1`
 - Дополнительные действия не требуются
 - Проверка `count` перед обновлением семафора
- Батч с `block.count > 1` без пролога
 - +1 атомарная инструкция на запуск
- Батч с `block.count > 1` с прологом
 - +2 атомарных инструкции на запуск первого батча

Scheduler v2.0 - результаты

- Интерфейс для **data-parallel** задач
 - `block.index` и `block.count`
 - Синхронизация (пролог и эпилог)
- Все еще низкий overhead
 - В случае не data parallel jobs – как и в v1.0
- Все еще простой код

- Многопоточный код бывает простым
 - Как код высокого уровня... (parallel_for)
 - ... так и код низкого уровня (MPMC FIFO)
- От job scheduler-а нужно немного
 - Простой интерфейс
 - Низкий overhead
- Будьте проще, и к вам потянутся такты.



Арсений “Zeux” Капулкин

CREAT Studios

a.kapoulkine@creatstudios.com