

Solution set

1. Reduce and Scan (**24 points**) One of the problems below can be solved using reduce, the other can be solved using scan. Identify which is which. For the problem that can be solved using reduce, describe the `Leaf` and `Combine`, and `Root` functions – the `Root` function may (or may not) be the identity function. For the problem that can be solved using scan, describe the `Leaf1` and `Leaf2`, and `Combine` functions.

- (a) (**12 points**) Given a list of numbers that is distributed across the worker processes, compute a new list whose i^{th} element is the maximum element up to and including the current one minus the minimum element up to and including this one. I'll call this function `gap` and here are a few examples:

```
gap([5,2,6,18,-3,4,8,-1,20]) -> [0,3,4,16,21,21,21,21,23];
gap([]) -> [];
gap([42]) -> [0];
```

Solution The problem is an instance of scan. In my solution, the `Leaf1` function reads the list and returns a tuple of the form `{Min, Max}` where `Min` is the minimum element of the list and `Max` is the maximum element. If the list is empty, `Leaf1` returns the atom '`empty`'. The `Combine` function takes two such tuples (or '`empty`' atoms) and returns the maximum and minimum for the combined subtrees (or '`empty`' if both are empty). The `Leaf2` function takes an `AccIn` value that is the tuple with the maximum and minimum for everything to the left of this leaf. It computes the gap for each element and stores it with the given key. I wrote a helper function, `gap(List, AccIn)` that returns a tuple of the form `{GapList, {Min, Max}}`, where `GapList` is the gap for each element of `List`, and the `{Min, Max}` tuple is the tuple produced by `Leaf1` or `Combine` (or the atom '`empty`').

```
gap(L, AccIn) -> mapfoldl(
  fun
    (E, empty) -> {0, {E, E}};
    (E, {Min, Max}) ->
      X = min(Min, E),
      Y = max(Max, E),
      {Y-X, {X, Y}}
  end, AccIn, L).

gap(W, Src, Dst) ->
  wtree:scan(W,
    fun(S) -> % Leaf1
      element(2, gap(workers:get(S, Src), empty))
    end,
    fun(S, AccIn) -> % Leaf2
      workers:put(S, element(1, gap(workers:get(S, Src), AccIn)))
    end,
    fun % Combine
      (empty, Right) -> Right;
      (Left, empty) -> Left;
      ({Lmin, Lmax}, {Rmin, Rmax}) ->
        {min(Lmin, Rmin), max(Lmax, Rmax)}
    end
  ).
```

Note: I haven't tried compiling and running this code yet – you didn't have that opportunity on the exam. When I do, there will almost certainly be some minor error. I'll fix them and might add

that version. You can think of the difference between the code written here and the final version as indicating the range of what I will consider to be an acceptable solution. Of course, there are many other valid implementations other than the one shown here. They will get full credit as well.

- (b) **(12 points)** Given a list of numbers, report the number of consecutive Pythagorean triples. If x , y , z are consecutive elements of the list, then they are a Pythagorean triple iff $x^2 + y^2 = z^2$. Here are few examples:

```
nPythag3([1,2,4,6,8,10,11,3,4,5,12,13,14]) -> 3; % [6,8,10], [3,4,5], [5,12,13]
nPythag3([1,-2,3,-4,5]) -> 1; nPythag3([0,1,0,1,0,1,0,1]) -> 3;
```

Solution The problem is an instance of reduce. In my solution, each leaf worker process produces a tuple of the form `{First2, N_triple, Last2}` where `First2` is a list of the first two elements of the worker's piece of the lists; `Last2` is a list of the last two elements of the worker's piece of the lists; and `N_triple` is the number of Pythagorean triples in the worker's piece of the list. If the worker's piece of the list has fewer than two elements, then `Leaf` just returns the list. The `Combine` function combines these in the “obvious” way. Have a different pattern for leaves with short lists makes handling the special cases fairly straightforward. The `Root` function extracts the count from the triple, or if the top-value is a list, returns 0. Finally, I wrote `nPythag3(List)`, a sequential implementation which is used by `Leaf` and `Combine`.

```
nPythag3(L=[A, B | _]) ->
{N, Last2} = nPythag3(L, 0),
{[A,B], N, Last2};
nPythag3(L) -> L.
nPythag3(W, Src) ->
wtree:reduce(W,
  fun(S) -> nPythag3(workers:get(S, Src)) end, % Leaf
  fun % Combine
    ({LL, LN, LR}, RL, RN, RR) ->
      {LL, LN + element(2, nPythag3(LR++RL)) + RN, RR};
    (L, {RL, RN, RR}) ->
      {First2, MidN, RL} = nPythag3(L++RL),
      {First2, MidN+RN, RR};
    ({LL, LN, LR}, R) ->
      {LR, MidN, Last2} = nPythag3(LR++R),
      {LL, LN+MidN, Last2};
    (L,R) -> nPythag3(L++R)
  end,
  fun % Root
    ({-, N, -}) -> N;
    (L) -> 0
  end
).
.
```

2. **Filter Locks (16 points)** I've mentioned several times that Peterson's mutual exclusion algorithm can be extended to apply to any number of clients. The algorithm is called a “filter lock”. Here's the code:

```
00: int level[N], victim[N]; // initially all 0
01: lock(i) {
02:   for(int j = 1; j < N; j++) {
03:     level[i] = j;
04:     victim[j] = i;
```

```

05:     while((victim[j] == i) && ∃ 0 ≤ k < N. (k ≠ i) && level[k] ≥ j);
06:   }
07: }

08: unlock(i) {
09:   level[i] = 0;
10: }

```

The key idea in the lock is that each thread must pass through $N - 1$ “levels” of locking (the `for(j...)` loop) before entering its critical section. To advance from level j to level $j+1$, there must be no other threads at level j or higher, **or** there must be some other thread that is the “victim” at level j . A consequence is that if there is another thread at this level or higher, then there must be at least one thread “left behind” at each of the lower levels. In particular, if thread x is in the critical section, and thread y is at level $N - 1$, there are threads at every level from 1 to $N - 1$; thread y is the victim at level $N - 1$; and thread y can’t advance to the critical section.

- (a) **(8 points)** Show that the filter lock is **not** starvation free. For example, consider a filter lock with three clients: 0, 1, and 2. Show that one of the clients can spin forever while the other two alternately acquire and release the lock.

Solution: Proof that filter locks are starvation free. If a thread is stuck spinning at line 05, then this thread must be the victim, and there must be another thread at or above this level. If there is another thread at this level, then that thread can advance because it is not the victim. This shows that the filter lock is deadlock free. Once the threads that are ahead of this one make it to the critical section and back to idle, this one must will be able to advance.

The concern is that some other threads may come reach this level in the meantime – for example, threads that were above this one could reach their critical section, come back around and then pass this one again. The key observation (and this is where I blew it) is that if thread i_1 is at level j , and thread i_2 enters level k then thread i_2 sets the `victim[j]` variable to i_2 . It will not get set back to i_2 (until after i_1 has reached its critical section and is trying again for the lock), because only thread i_1 can set a victim variable to i_1 . Thus, once some other thread enters level j , thread i_1 can progress at least to level $j + 1$. Eventually, thread i_1 gets the lock.

- (b) **(8 points)** Show that if the statements at lines 03 and 04 are exchanged (i.e. `victim[j] = i;` is performed *before* `level[i] = j;`, then mutual exclusion can be violated.

Solution: The following execution shows a violation:

0. All threads are at level 0.
1. Thread 1 sets `victim[1] = 1;`
2. Thread 2 sets `victim[1] = 2;`
3. Thread 2 sets `level[2] = 1;`
4. Thread 2 checks the condition at line 05. It is the victim, but there is no other thread with `level` at or above 1.
5. Thread 2 enters its critical section.
6. Thread 1 sets `level[1] = 1;`
7. Thread 1 checks the condition at line 05. It is the not victim.
8. Thread 1 enters its critical section.

3. Sorting (25 points)

The question finishes showing that the compare-and-swap elements of sorting networks can be replaced by “merge-and-split” operations. Homework 4, question 2, also looked at this. This can be used to obtain practical parallel sorting algorithms for common parallel architectures.

Let M be a positive integer. We will consider merge-and-split operations on lists of length M . In particular,

```
merge_and_split({In0, In1}) -> {Out0, Out1}
```

where:

In0 and In1 are lists of M elements in ascending order.

Out0 is the first M elements of merge(In0, In1), sorted into ascending order.

Out1 is the last M elements of merge(In0, In1), sorted into ascending order.

For example:

```
merge_and_split({[1, 4, 9, 16, 25, 36], [3, 5, 7, 9, 11, 13]}) ->
{[1, 3, 4, 5, 7, 9], [9, 11, 13, 16, 25, 36]}.
```

- (a) **(5 points)** What is the result of

```
merge_and_split({[0, 0, 0, 0, 1, 1], [0, 0, 0, 1, 1, 1]})
```

?

Solution: {[0, 0, 0, 0, 0], [0, 1, 1, 1, 1, 1]}.

- (b) **(5 points)** Show that if In0 and In1 are lists of 0s and 1s in ascending order, and

```
{Out0, Out1} = merge_and_split({In0, In1})
```

then at least one of Out0 or Out1 must be clean (i.e. all 0s or all 1s).

Solution: There are a total of $2M$ elements in In0 and In1. If at least M of these are 0s, then Out0 is all 0s and therefore clean. Otherwise, at least $M + 1$ of the input elements are 1s, and Out1 is all 1s and therefore clean.

- (c) **(10 points)** Let S be a sorting network with N inputs and N outputs. Let S' be the network obtained by replacing every compare-and-swap in S with a merge-and-split operation. Now, let A_0, A_1, \dots, A_{N-1} be any N sorted lists of 0s and 1s of length M . Apply these lists as the inputs to S' , and let X_0, X_1, \dots, X_{N-1} be the resulting outputs.

Show that there is an input to S' , B_0, B_1, \dots, B_{N-1} such that B is a permutation of X and when B is applied as the input to S' , X is the output.

Solution: The main idea is that we can derive B by working backwards through the sorting network from X . We use the forward computation that started with A to see where the clean lists are in the network. If A produces a clean list at the output of some merge-and-split operation, we'll keep that list in our backwards pass for deriving B . We note that if A produces a dirty list at the output of some merge-and-split operation, we can change the number of 1s or 0s in that list, and still be able to find a valid input to that merge-and-split that produces the outputs that we are using in our backwards sweep. This lets us keep the lists as a permutation of the lists from X , and we never have to break a clean lists from X into the dirty lists of A that went into it. Quite a bit of partial credit will be given for making observations like those above.

The rest of my solution gives the details for the proof. I give lots more details than I require in a full-credit solution. I'm trying to write a solution that will be clear to someone who has already thought about the problem but may not have solved it.

We can describe a sorting network as a list of tuples. The list $\{\{I, J\} \mid T1\}$ describes a network that first does a merge-and-split of lists I and J , and applies the resulting lists to the sorting network described by $T1$. In particular, Let

```
derive_b(A, []) -> A;
derive_b(A, S1 = [{I, J} | T1]) ->
  {AI, AJ} = {nth(I, A), nth(J, A)},
  {RI, RJ} = merge_and_split({AI, AJ}),
  R = replace(I, RI, replace(J, RJ, A)),
  Y = derive_b(R, T1),
```

```

{YI, YJ} = {nth(I, Y), nth(J, Y)}
{BI, BJ} = if
    clean(AI) and (YJ == AI) -> YJ, YI;
    clean(AJ) and (YI == AJ) -> YJ, YI;
    true -> YI, YJ
end,
B = replace(I, BI, replace(J, BJ, Y)).

```

Where `nth(I, List)` is the I^{th} element of `List`; `replace(I, X, List)` is the list obtained by replacing the I^{th} element of `List` with `X`; and `clean(List)` is true if every element of `List` is 0 or every element of `List` is 1. If `S` is a sorting network (using merge-and-split operations) and `LL` is a list of lists, then I'll write `S(LL)` to denote the list of lists obtained by applying sorting network `S` with input `LL`.

To understand the code, note that $\{I, J\}$ are the indices of the lists for the next merge-and-split. Thus `AI` and `AJ` are the inputs to that merge-and-split, and `RI` and `RJ` are the outputs. The goal now is to derive `B` so that $S_1(A) = S_1(B)$ and the lists of `B` are a permutation of the lists of $S_1(A)$. The code does this recursively. It constructs `Y` so that $T_1(Y) = S_1(A)$ and the lists of `Y` are a permutation of the lists of $S_1(A)$. Then, it derives `B` from `A` and `Y`.

Let `B = derive_b(A, S')`. I'll show that `B` has the properties described above. My proof is by induction on S_1 , and my induction hypothesis is

$S_1(B) == S_1(A)$,

and The lists of `B` are the same (but possibly permuted) as the lists of $S_1(A)$,

and If the I^{th} list of `A` is clean, then the I^{th} list of `B` is the same as the I^{th} list of `A`.

The base case is when $S_1 = []$, then $B = A = X$, and the claims are immediate. Otherwise, let $S_1 = [\{I, J\} \mid T_1]$. From the induction hypothesis, we have that $T_1(R) = S_1(A)$. There are two cases to consider:

`clean(AI) and clean(AJ)`: then `RI` and `RJ` are both clean. From the induction hypothesis, $YI=RI$ and $YJ=RJ$ are both clean. From the declarations for `BI` and `BJ`, $BI=AI$ and $BJ=AJ$, therefore $B=A$ and $S_1(B) = S_1(A)$, and all clean lists of `A` are equal to their counterparts in `B`. From the declarations for `BI` and `BJ`, `B` is a permutation of `Y`. From the induction hypothesis, `Y` is a permutation of $S_1(A)$. Therefore, `B` is a permutation of $S_1(S)$.

Otherwise: One of `AI` or `AJ` (or both) are dirty. From question 3b, one of `RI` or `RJ` is clean. If `AI` is clean and all 0s, then `RI` will be clean and all 0s (def. `merge_and_split`); `YI` will be clean and all 0s (induction hypothesis); and `BI` will be clean and all 0s (declaration of `BI`). Similar arguments apply in the other three cases where `AI` or `AJ` are clean. If `AI` and `AJ` are both dirty, then one of `RI` or `RJ` will be clean (question 3b); and the corresponding `YI` or `YJ` will be clean. The other one may be either clean or dirty. In this case $BI=YI$ and $BJ=YJ$. By construction,

$\{YI, YJ\} = \text{merge_and_split}(\{BI, BJ\})$

Therefore, $S_1(B) = T_1(Y) = S_1(A)$. Furthermore, the lists of `B` are a permutation of the lists of $S_1(A)$ by the same argument as used in the previous case.

This completes the induction proof. Therefore, `B = derive_b(A, S')` satisfies the induction hypothesis, which means that it is a permutation of `X = S'(A)` and `S'(B) = X` as required.

- (d) **(5 points)** The induction argument from question 3c shows that if the lists B_1, \dots, B_{N-1} are applied as the inputs to S' , then for each `merge_and_split` operation in S' , either $\{Out0, Out1\} = \{In0, In1\}$, or $\{Out0, Out1\} = \{In1, In0\}$. For $0 \leq i < N$, let C_i be the number of 1s in list B_i . Show that if S' sorts the B lists incorrectly, S must sort C incorrectly as well.

Solution: If `B` is applied as an input to `S'` then each merge-and-split operation has at least one clean input – otherwise, the number of 1s in the lists wouldn't be preserved by the merge-and-split operations. If one input to a merge-and-split is clean, then it must have a clean output, and the other output must match the other input. In other words, the merge and split either copies $\{In0, In1\}$ to $\{Out0, Out1\}$

`Out1}` or it swaps them. In the problem statement, I told you you could assume this; so, a solution can get full credit without the argument of this paragraph.

The lists consist of 0s and 1s and are sorted into ascending order; so, there are only $M + 1$ possible lists, depending on how many 1s are in the list. Thus, each input to a merge-and-split can be described using integers in $0, \dots, M$, and the outputs of the merge-and-split are the same of the inputs if `In0` has fewer 1s than `In1` and swapped if `In0` has more 1s than `In1`. Thus, a compare-and-swap of the integers is equivalent to merge-and-split of the lists. This can be applied to the entire sorting network. Therefore, if the merge-and-split sorting network produces an unsorted result with the B list-of-lists, then the compare-and-swap network will produce an unsorted result with the C list-of-integers.

Hint: you can do this even if you didn't solve question 3c. In particular, you are allowed to assume what I said about the `merge_and_split` operations.

4. Matrix Multiplication (20 points)

Let's say you get a job writing code at a local company that needs to do some moderately intense number crunching. A typical problem is computing products of $N \times N$ matrices, where N ranges from 1000 to 20000.

- (a) **(4 points)** Let's say you've got a laptop with a 2.5GHz, quad-core processor. From mini-assignment 4, we know that multiplying a $N_1 \times N_2$ by a $N_2 \times N_3$ matrix on a machine with a clock frequency of f takes time of roughly $3N_1 N_2 N_3 / f$. What is the time to compute a matrix products for $N = 1000$, $N = 5000$, and $N = 20000$?

Solution: The time is $3N^3/f$. For $f = 2.5$ GHz, this is $(1.2 \text{ seconds})(N/1000)^3$.

$$N = 1000: T_{seq} = 1.2 \text{ seconds.}$$

$$N = 5000: T_{seq} = 1.2 * (5^3) \text{ seconds} = 1.2 * 125 \text{ seconds} = 150 \text{ seconds.}$$

$$N = 20000: T_{seq} = 1.2 * (20^3) \text{ seconds} = 1.2 * 8000 \text{ seconds} = 9600 \text{ seconds.}$$

- (b) **(4 points)** You try to convince your manager to buy a cluster of linux machines for your computation. Your proposal is for 32, 8-core processors running at 2.5GHz. The total number of processors is 256. For simplicity, we'll pretend that there isn't any multi-threading. Assume that the time to transfer a block m double-precision values between processors (over the network, it's a cluster) is:

$$T_{network} = 0.1\text{ms} + m * 50\text{ns}$$

where $1\text{ms} = 10^{-3}\text{seconds}$, and $1\text{ns} = 10^{-9}\text{seconds}$.

Using the algorithm we derived in class (Sept. 24), each processor performs P matrix products where it multiplies a $(N/P) \times (N/P)$ matrix by a $(N/P) \times N$ matrix. Furthermore, each processor sends (and receives) $P - 1$ messages of size N^2/P double precision numbers. For simplicity, assume that the time for sending messages cannot be overlapped with computation time. What is the speed-up when computing matrix products for $N = 1000$, $N = 5000$, and $N = 20000$?

Solution: Each processor computes P multiplications of a $(N/P) \times (N/P)$ matrix by a $(N/P) \times N$ matrix. Each matrix multiplication takes $3N^3/(P^2f)$ time. The total compute time for each processor is $P \frac{3N^3}{P^2f} = 3N^3/(Pf)$ which is just T_{seq}/P . The communication time is $(P - 1) * (0.1 \text{ ms} + 50(N^2/P) \text{ ns})$. That gives us the following table:

N	T_{seq}	$T_{par,1}$	$T_{par,2}$	$T_{par,3}$	T_{par}	speed-up
1000	1.2s	4.6875ms	1.5ms	49.8ms	55.9875ms	21.4
5000	150s	586ms	1.5ms	1.245s	1.8326s	81.9
20000	9600s	37.5s	1.5ms	19.92s	57.42s	167.2

where $T_{par,1} = 3N^3/(Pf)$ is the compute time for the parallel version; $T_{par,2} = (P - 1) * 0.1 \text{ ms}$ is the part of the communication time that is independent of matrix size; and

$$T_{par,3} = (P - 1)50 \frac{N^2}{P} \text{ ns} = 50 \frac{P-1}{P} \left(\frac{N}{1000} \right)^2 \text{ ms}$$

is the part of the communication time that depends on the matrix size. Of course, to save time, you didn't need to copy numbers from your solution to question 4a. I put them all into the table to make it easier to read.

Note: for all parts of this problem, I did most of the calculations in my head, and used a calculator for the multi-digit divisions. My answers may be slightly off by one or two in the least significant digits. That's indicative of the error-tolerance that I'll include in grading.

- (c) **(4 points)** Your manager isn't convinced, but you really want to her to buy that linux cluster. You read up on algorithms for matrix multiplication and find an algorithm where each processor performs \sqrt{P} matrix products where it multiplies a $(N/\sqrt{P}) \times (N/\sqrt{P})$ matrix by another $(N/\sqrt{P}) \times (N/\sqrt{P})$. Furthermore, each processors sends (and receive) \sqrt{P} messages of size N^2/P double precision numbers. With the new algorithm, what is the speed-up when computing matrix products for $N = 1000$, $N = 5000$, and $N = 20000$?

Solution: The computation time stays the same. There are now \sqrt{P} messages sent per processor instead of $P - 1$. I'll just scale the communication times ($T_{par,2} + T_{par,3}$) from the previous problem by $\sqrt{P}/(P - 1) = 0.62745$. The new table is:

N	T_{seq}	$T_{par,1}$	$T_{par,2+3}$	T_{par}	speed_up
1000	1.2s	4.6875ms	3.22ms	7.91ms	151.7
5000	150s	586ms	78.2ms	664.2ms	225.8
20000	9600s	37.5s	1.25s	38.75s	247.7

Reducing the number of messages makes a big difference!

- (d) **(8 points)** A coworker installs an optimized BLAS (numerical library) which is faster than the simple code I showed in class. You find that the time to multiply a $N_1 \times N_2$ matrix by a $N_2 \times N_3$ matrix on a machine with a clock frequency of f now takes time of roughly $1.2N_1N_2N_3/f$. What is the sequential time (as in question 4a) and parallel time (as in question 4c when the BLAS library is used)?

- What is the sequential time (as in question 4a) and parallel time (as in question 4c when the BLAS library is used)?

Solution: The sequential times get divided by $3/1.2 = 2.5$. For the parallel version, the compute time is divided by 2.5, but the communication time stays the same.

N	T_{seq}	$T_{par,1}$	$T_{par,2+3}$	T_{par}	speed_up
1000	0.48s	1.875ms	3.22ms	4.1ms	117.5
5000	60s	234.4ms	78.2ms	312.6ms	191.9
20000	3840s	15s	1.25s	16.25s	236.3

- Write one sentence to summarize the impact of the faster sequential algorithm on the execution time.

Solution: Execution times dropped for both the sequential and parallel versions.

- Write one sentence to summarize the impact of the faster sequential algorithm on the speed-up.

Solution: Speed-ups decreased as well.

- Write one sentence to explain the relationship you observe between execution time and speed-up?

Solution: A faster sequential implementation results in lower speed-ups because the overheads (communication) become a larger percentage of the total execution time for the parallel program.

5. Other Stuff (16 points)

Answer each question below with 1-3 sentences. Points may be taken off for long answers.

- (a) **(4 points)** What is instruction level parallelism?

Solution: Executing multiple instructions in parallel by identifying their dependencies and executing instructions when their dependencies are resolved. Super-scalar architectures are one way to exploit instruction level parallelism.

(b) **(4 points)** What is “single-instruction, multiple-thread” parallelism?

Solution: Fetching and decoding a single instruction stream and sending it to multiple execution pipelines. Each pipeline can resolve branch conditions independently. Instead of taking or not taking a branch, instructions are executed conditionally based on the branch outcome. This is called “predicated execution.” GPUs are an example of an architecture that provides single-instruction multiple-thread parallelism.

(c) **(4 points)** How does map-reduce handle machine and network failures?

Solution: The map and reduce processes can be (re-)executed on a different machine if the original machine fails to respond (e.g. to a ping). The entire job can be re-executed with a different master process if the master fails to respond.

(d) **(4 points)** What is a memory fence?

Solution: An operation that forces pending memory operations to complete, thus ensuring that all cache coherence operations have finished as well.

Note, my answer for single-instruction, multiple-thread was five sentences. A three sentence answer is fine as long as it clearly gets two or three key points. The same applies for the other problems as well.

This exam has 10 pages (plus the cover sheet). Please check that you have a complete exam.

Answer question 0, and any **four** out of questions 1–5. If you write solutions for all questions, please indicate which you want to have graded; otherwise the graders will select which one to discard. Note that questions 1–3 are worth 20 points each and questions 4–5 are worth 24 points each. Depending on which questions you choose to answer, your maximum score will be either 86 or 90; however, no matter which questions you choose to answer the exam will be graded on a scale of 86 points.

0. (2 points)

Sign below to confirm that you have read and understood the exam instructions.

I, Mark Greenstreet, confirm that I have read and understood the exam instructions and that I only need to attempt **four** questions out of questions 1–5 below. If I attempt all questions, I understand that I should clearly indicate which ones to grade. If I attempt all questions and do not indicate which ones to grade, then I accept that the graders will choose at their own convenience which to grade.

1. **Erlang** (20 points)

The final two pages of the exam contain Erlang definitions of the functions `last`, `lastv2` and `biggest_product`, along with associated functions that are used to define these functions. You may find it convenient to tear off the last sheet of the exam so you can refer to these functions while answering this question.

For each of these functions, state its asymptotic worst-case run time (using big-O notation) and briefly justify your claim (you do not need to provide a formal proof or recurrence).

For your run time analysis you may use the variable N to refer to the length of the list, the variable P to refer to the size of the worker pool \mathbb{W} , and variable λ to refer to the ratio between the time taken to complete a global action (such as communicating a message between workers) and the time taken to complete a local action (such as a basic arithmetic operation). You may assume that the following operations are all unit time local actions:

- `hd`, `tl`, pattern matching
- `+`, `-`, `*`, `/`, `abs`, `div`, `max`, `rem`
- `<`, `=<`, `=:=`, `==`, `/=`, `=/=`, `>=`, `>`
- `and`, `andalso`, `or`, `orelse`, `not`.

You may also assume that `lists:foldl(Fun, Acc0, List)` takes time linear in the length of `List` plus the time for the (linear number of) calls to `Fun`.

The possible run-time complexities are: $\mathcal{O}(\log N)$, $\mathcal{O}(N)$, $\mathcal{O}(N \log N)$, $\mathcal{O}(N^2)$, $\mathcal{O}(N + P)$, $\mathcal{O}(N + \lambda)$, $\mathcal{O}(NP)$, $\mathcal{O}(\lambda N)$, $\mathcal{O}((N/P) + \log P)$, $\mathcal{O}((N/P) + \lambda \log P)$, $\mathcal{O}((N/P) + \lambda)$, or $\mathcal{O}((N/P) + \lambda P)$. Not all of these complexities occur in this question.

(a) (5 points) `last`: return the last element of a list. Asymptotic complexity and brief justification:

$\mathcal{O}(N)$. Each recursive call to `last` takes $\mathcal{O}(1)$ time. There are N such recursive calls.

Note: from the review problems.

(b) (5 points) `lastv2`: another implementation of `last`. Asymptotic complexity and brief justification:

$\mathcal{O}(N^2)$. Each recursive call to `lastv2` takes $\mathcal{O}(\text{length(List)})$ to check the “`when length(List) > 1`” part of the guard. This is done for `length(List) = N, N-1, ..., 1`. The total time is

$$\sum_{i=1}^N i = \frac{N^2 + N}{2} \in \mathcal{O}(N^2)$$

Note: from the review problems.

- (c) (10 points) `greatest_product`: return the greatest product of pairs of adjacent elements in a list (parallel implementation). Asymptotic complexity and brief justification:

$$O\left(\frac{N}{P} + \lambda \log P\right).$$

Each worker calls `bp_leaf` with a list of $\frac{N}{P}$ elements. `bp_leaf` does $\mathcal{O}(1)$ work plus the work of `gp(List)`. `gp(List)` calls a tail recursive version of itself. The tail-recursive `gp(List, Acc)` does $\mathcal{O}(1)$ work per recursive call – note that `my_max` takes $\mathcal{O}(1)$ time. Thus `gp(List, Acc)` and therefore `gp(List)` and `bp_leaf(List)` take $\mathcal{O}(\text{length(List)}) = \mathcal{O}(N/P)$ time.

Each level of the combine tree involves one send-and-receive pair at a cost of λ and $\mathcal{O}(1)$ sequential computation. There are $\lceil \log_2 P \rceil$ levels to the tree. The total cost for the combine is $\mathcal{O}(\lambda \log P)$.

`pb_root` extracts the second element from a three element tuple. This takes $\mathcal{O}(1)$ time.

Combining the results for `bp_leaf`, `bp_combine`, and `bp_root` yields the total time of $O\left(\frac{N}{P} + \lambda \log P\right)$.

Note: We've covered this in class many times. Just writing the answer with a reference to the lecture slides is a sufficient justification. For example, "See Jan. 13 lecture slides, slide 3".

2. `largest_gap(W, Key)` (20 points)

`largest_gap` should return the largest gap between adjacent numbers of the list of numbers stored on the worker processes of `W` and associated with `Key`. Here is the sequential version:

```
largest_gap([]) -> undef;
largest_gap([H | Tail]) ->
  {Gap, _} = lists:foldl(fun(X, {G, PrevX}) -> {my_max(abs(X-PrevX), G), X} end,
                        {undef, H}, Tail),
  Gap.
```

where `my_max` is the function defined in question 1. For example,

```
largest_gap([]) -> undef.
largest_gap([1,4,9,16,20]) -> 7.  % |16 - 9| = 7
largest_gap([1,4,9,16,2]) -> 14.  % |2 - 16| = 14
```

The parallel version of `largest_gap` is a typical reduce pattern, and you will write the functions `lg_leaf`, `lg_combine` and `lg_root` in the following implementation:

```
largest_gap(W, Key) ->
  wtree:reduce(W,
    fun(ProcState) -> lg_leaf(wtree:get(ProcState, Key)) end,
    fun(Left, Right) -> lg_combine(Left, Right) end,
    fun(Top) -> lg_root(Top) end
  ).
```

Assuming that there are N elements in the list spread evenly across P workers in the worker pool `W` (where $P \ll N$), your implementation should achieve a speedup close to P . Your implementation should fail if there is no value associated with `Key` or if the value associated with `Key` is not a list of numbers.

You may call any of the functions defined in question 1, the sequential `largest_gap(List)` function defined above, any Erlang built-in functions, and any functions from the `lists`, `misc`, `workers`, or `wtree` modules. **Hint:** Some functions from question 1 may prove inspirational even if they cannot be called directly.

- (a) (8 points) Your code for `lg_combine(Left, Right)`:

```
% All three parts of this problem can be solved by copying the code
% for biggest_product from Q1 with a few small changes.
lg_combine(empty, Right) -> Right; % adapted from bp_combine
lg_combine(Left, Empty) -> Left;
lg_combine({LL, LBig, LR}, {RL, RBig, RR}) ->
  {LL, my_max([LBig, abs(LR-RL), RBig]), RR}. % my_max from Q1.
```

(b) (8 points) Your code for `lg_leaf(Data)`:

```
% adapted from bp_leaf.
% largest_gap(List) is from the problem statement.
lg_leaf([]) -> empty;
lg_leaf(List) -> {hd(List), largest_gap(List), tl(List)}.
```

(c) (4 points) Your code for `lg_root(Top)`:

```
lg_root(V) -> bp_root(V).
```

3. Message Passing Network Topologies (20 points)

We are building a message passing machine with 64 nodes, and we need to decide how to connect its network. For each of the following topologies, specify the bisection bandwidth, diameter and maximum number of network ports needed by a node.

(a) (3 points) Ring.

bisection bandwidth: 2

diameter: 32

network ports: 2

(b) (3 points) 2D mesh.

bisection bandwidth: 8

Note: \sqrt{P} assuming a simple mesh, not toroidal because the problem didn't use the word "toroidal".

Draw a line that separates the left and right halves of the mesh.

diameter: 14

Note: 7 hops in each dimension.

network ports: 4

Note: connections to north, south, east, and west neighbours.

(c) (3 points) 3D mesh.

bisection bandwidth: 16

Note: $\sqrt[3]{P}$ assuming a simple mesh, not toroidal because the problem didn't use the word "toroidal".

Draw a plan that separates the top half of the mesh from the bottom half.

diameter: 9

Note: 3 hops in each dimension.

network ports: 6

Note: same connections as above, plus up and down.

(d) (3 points) 6D hypercube.

bisection bandwidth: 32

diameter: 6

Note: example, node 0 to node 63.

network ports: 6

Note: one for each bit of the node index.

(e) (4 points) State **two** advantages of the 6D hypercube over the 2D mesh? Answers which give more than two will receive zero.

- Small diameter → low routing latency.
- High bisection → fewer communication bottlenecks.
- Note: for the solution set, I'll include another response that gets full credit:
 - No special cases at edges.

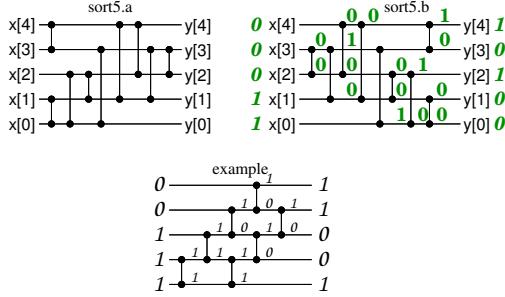


Figure 1: Sorting Networks

- (f) (4 points) State **two** advantages of the 2D mesh over the 6D hypercube? Answers which give more than two will receive zero.

- Easily mapped to 2D technologies such as chips or circuit boards.
- Bounded degree nodes make the design more modular.
- Note: for the solution set, I'll include another response that gets full credit:
 - Doesn't become "all wire" for networks with a large number of processors.

4. Sorting Networks (24 points)

Figure 1 shows two sorting networks with five inputs: sort5.a, and sort5.b. One of these networks sorts correctly and one does not.

- (a) (10 points) Identify which of the networks **does not** sort correctly. Label the inputs with values of **0** or **1** that the network does not sort correctly. Label the output of each compare-and-swap with its value for this input, and label the outputs of the sorting network with their values. The "example" network provides an example of such a labeling (and also an incorrect attempt at sorting).

Which network do you claim is incorrect (circle one):

sort5.a

sort5.b

Note: The counter-example has 1s for $x[0]$ and $x[1]$ and 0s for $x[2]$, $x[3]$ and $x[4]$. The compare-and-swap (1,3) moves the 1 for $x[1]$ into row 3. Note that the only path for the 1 for $x[0]$ to rows 3 or 4 is through the compare-and-swap (0,3). However, it is blocked by the 1 from $x[1]$ that has already moved to row 3. This shows that the network cannot sort this input correctly.

Remember to label that network with a counter-example in the figure above.

- (b) (14 points) For the other network, explain specifically how it sorts any combination of three **0**s and two **1**s correctly.

We observe that the three compare-and-swap modules in the bottom left, i.e. (0,1), (0,2), and (1,2) sort rows 0, 1, and 2 into ascending order. Likewise, the compare-and-swap module in the top-left, i.e. (3,4), sorts rows 3 and 4 into ascending order. Call these four compare-and-swap operations "the initial compare and swaps"

Now, consider the number of 1s in $\{x[3], x[4]\}$.

zero 1s: In this case, the values of rows 0 to 4 after the initial compare and swaps must be (0,1,1,0,0). The 1 in row 2 moves to row 4 with the compare-and-swap (2,4), and the 1 in row 1 moves to row 3 with the compare-and-swap (1, 3). Thus, the two ones make it to the top two rows.

one 1: In this case, the values of rows 0 to 4 after the initial compare and swaps must be (0,0,1,0,1). The 1 in row 2 moves to row 3 with the final compare-and-swap (2,3).

two 1s: In this case, the values of rows 0 to 4 after the initial compare and swaps must be (0,0,0,1,1). The data are already sorted, and all of the compare-and-swaps preserve this order.

Compare-and-swap operations preserve the total number of 1s in the network, and the total number of 0s. Showing that the two 1s arrive at $y[3]$ and $y[4]$ implies that $y[0] = y[1] = y[2] = 0$.

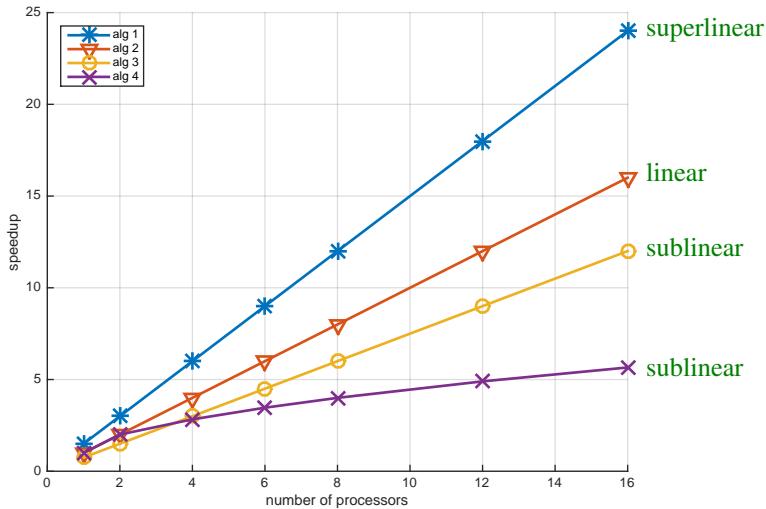


Figure 2: Speedup for some parallel algorithms.

5. Short Answer (24 points)

- (a) (7 points) Figure 2 shows the speedup measured for four algorithms: alg 1, alg 2, alg 3 and alg 4. In the space to the right of the plot, label each curve as either *sublinear*, *linear* or *superlinear* speedup. Not all speedups may be present in the plot.

Give **two** general reasons why parallel programs sometimes display sublinear speedup *other than computational or communication overhead*. Answers which give more than two reasons will receive zero.

- Non-parallelizable code.
- Extra computation.
- Note: for the solution set, I'll include other responses that get full credit:
 - memory overhead
 - synchronization overhead
 - resource contention
 - idle processors

Give **one** general reason why parallel programs sometimes display superlinear speedup. Answers which give more than one reason will receive zero.

- more fast memory (in total)
- Note: the “fast memory” answer is the right explanation 99% of the time. occasionally, a problem is so naturally parallel that the sequential algorithm has control overhead to serialize the operations that the parallel code avoids. A common explanation is “invalid comparison” because there’s something wrong with the way the times were measured, but I don’t think that’s an intended answer.

- (b) (6 points) For sequential computing there is the successful and near ubiquitous Von Neumann architecture; in the early days of computing there were other sequential architectures, but they have essentially disappeared from use. In contrast, we have discussed three very different approaches to achieving parallel computation that are still widely used, often in combination with one another: superscalar, shared memory and message passing. For **each** of the three, give **one** reason why it has **not** become dominant over the others; in other words, one significant disadvantage. Note that your reasons need not be unique. Answers which give more than one reason for any approach will have marks deducted.

superscalar: Limited instruction level parallelism.

shared memory: Poor scaling of coherence models to large numbers of processors.

message passing: Network delays.

- (c) (6 points) You should remember from your algorithms course that sequential merge sort requires $\mathcal{O}(N \log N)$ time and $\mathcal{O}(N \log N)$ comparisons for N elements. We saw that under reasonable assumptions a bitonic sorting network requires $\mathcal{O}((\frac{N}{P})(\log N + \log^2 P))$ time and $\mathcal{O}(N \log N \log^2 P)$ comparisons for N elements with P processors. Assume now that we have M (where $M \gg P$) independent sorting problems each of size N . In this question we consider two methods of solving these M problems: either connecting the P processors together into a single bitonic sorting network and running the problems through the network one at a time, or running up to P independent sequential merge sorts at a time. We will assume in each case that the data is already distributed in a manner which is optimal for that case. Is the *latency* of a bitonic sorting network better, worse or about the same as independent sequential merge sorts? Is the *throughput* of a bitonic sorting network better, worse or about the same as independent sequential merge sorts? Briefly explain your answers.

Bitonic latency (circle one and briefly explain): BETTER WORSE SAME

The latency for parallel bitonic sort is $\mathcal{O}(\frac{N}{P} \log N \log^2 P)$ which is less than the latency for sequential merge sort, $\mathcal{O}(N \log N)$ when $N \gg P$.

Note: I found this one a little confusing. The problem statement says

“Assume now that we have M (where $M \gg P$) independent sorting problems each of size N .”

Is latency the time to solve one of these problems are all M of them? I can suggest rephrasings of the problem that would make this unambiguous, but it’s a bit late for that.

Bitonic throughput (circle one and briefly explain): BETTER WORSE SAME

With $M \gg P$, we can keep P processors busy working on independent merge sort problems. Thus, the bottleneck is the execution time of the processors. Sequential merge sort requires fewer total processor cycles than bitonic sort *and* it avoids the communication costs of the parallel algorithm. Thus, sequential merge sort will achieve higher throughput.

- (d) (5 points) We developed a tree-based reduce operation in Erlang—a language which assumes a message passing architecture, although it can be run on either message passing or shared memory hardware—to speed up parallel computations that produce a small amount of summary information at a single process from a large amount of data spread across many processes. If we were to use a language and hardware which supported shared memory, is a tree-based reduce still useful to speed up parallel computation? Briefly explain your answer. (If it is relevant to your answer, you may assume an idealized memory where access to any item costs the same amount for any process.)

The tree based algorithm is still faster. We need some kind of synchronization between the threads, and that comes at a cost of λ per synchronization. A tree has a cost of $\lambda \log P$ for synchronization. If a single master processor synchronized with each worker, the synchronization cost would be λP .

Note: I’m a bit concerned about “idealized memory”. If I have atomic reads and writes that take $\mathcal{O}(1)$ time, then I can implement a mutual-exclusion algorithm (such as Peterson’s) that provides lock acquisition or transfer in $\mathcal{O}(1)$ time. This removes the λ s from the analysis. In this case, a

tree-algorithm takes time $\mathcal{O}(\log P)$ and a singel master process takes tiem $\mathcal{O}(P)$. Again, the tree is faster, but it is likely that P is much less than the time for the worker processes to do their tasks before meeting at a reduce.

Functions for Q1a and Q1b

Do not write your answer on this page—it will not be graded. You may find it convenient to tear this page off when answering questions 1 and 2. If you tear it off, you need not submit it with the rest of your exam.

```
%%%%% Q1.a
% last(List) -> LastElement.
%   Return the last element of List. If List is not
%   a list or if List is empty, then an error will be thrown.
% Example:
%   last([1, 4, 9, 16]) -> 16.
%   For run-time analysis, N = length(List).
last([E]) -> E;
last([_ | T]) -> last(T).
```

```
%%%%% Q1.b
% lastv2(List) -> LastElement.
%   Another implementation of last(List).
%   For run-time analysis, N = length(List).
lastv2(List) when length(List) > 1 -> lastv2(tl(List));
lastv2([X]) -> X.
```

Repeated from the question 1 problem statement for your convenience:

For your run time analysis you may use the variable N to refer to the length of the list, the variable P to refer to the size of the worker pool W , and variable λ to refer to the ratio between the time taken to complete a global action (such as communicating a message between workers) and the time taken to complete a local action (such as a basic arithmetic operation). You may assume that the following operations are all unit time local actions:

- `hd`, `tl`, pattern matching
- `+, -, *, /, abs, div, max, rem`
- `<, =<, =:=, ==, /=, =/=, >, >=, >`
- `and, andalso, or, orelse, not.`

You may also assume that `lists:foldl(Fun, Acc0, List)` takes time linear in the length of `List` plus the time for the (linear number of) calls to `Fun`.

The possible run-time complexities are: $O(\log N)$, $O(N)$, $O(N \log N)$, $O(N^2)$, $O(N+P)$, $O(N+\lambda)$, $O(NP)$, $O(\lambda N)$, $O((N/P) + \log P)$, $O((N/P) + \lambda \log P)$, $O((N/P) + \lambda)$, or $O((N/P) + \lambda P)$. Not all of these complexities occur in this question.

Function for Q1c

Do not write your answer on this page—it will not be graded. You may find it convenient to tear this page off when answering questions 1 and 2. If you tear it off, you need not submit it with the rest of your exam.

```
%%% Q1.d
% greatest_product(W, Key) -> Number.
% W is a worker-tree. Key is the key for a list of numbers distributed
% across the workers. We return the largest product of adjacent numbers.
% If the list is empty or a singleton, then greatest_product returns 'undef'.
% If there is no value associated with Key, or if the value associated with
% Key is not list of numbers, then greatest_product fails.
% Example:
% If Key corresponds to the list [1, 2, 7, 3, 5, 4, 2, -8, 6, 3],
% Then the products of adjacent numbers are
% [2, 14, 21, 15, 20, 8, -16, -48, 18],
% and the largest product of adjacent numbers is 21.
% For run-time analysis: let N denote the total length of the list
% corresponding of Key; let P denote the number of processes.
% Assume that the list is divided into equal sized segments.
% Assume that there are enough processors so that all processes
% can run at the same time.
greatest_product(W, Key) ->
    wtree:reduce(W,
        fun(ProcState) -> bp_leaf(wtree:get(ProcState, Key)) end,
        fun(Left, Right) -> bp_combine(Left, Right) end,
        fun(Top) -> bp_root(Top) end
    ) .

bp_leaf([]) -> empty;
bp_leaf(List) -> {hd(List), gp(List), last(List)}.

bp_combine(empty, Right) -> Right;
bp_combine(Left, empty) -> Left;
bp_combine({LL, LBig, LR}, {RL, RBig, RR}) ->
    {LL, my_max([LBig, LR*RL, RBig]), RR}.

bp_root({_, Big, _}) when is_number(Big) -> Big;
bp_root(_) -> undef.

% gp(List) -> Number.
% a sequential version of greatest_product.
gp(List) -> gp(List, undef).
gp([], BigProd) -> BigProd;
gp([_], BigProd) -> BigProd;
gp([A | Tail = [B | _]], BigProd) ->
    gp(Tail, my_max(A*B, BigProd)).

% my_max: calculate max where 'undef' is less than any number
my_max(undefined, Y) -> Y;
my_max(X, undefined) -> X;
my_max(X, Y) -> max(X, Y).

% my_max for lists.
my_max(List) -> lists:foldl(fun(X, Acc) -> my_max(X, Acc) end, undefined, List).
```

1. The zero-one principle (**20 points + 5 points Extra Credit**)

Let $X[0..N - 1]$ be a *bitonic* array of N elements, where N is even, and every element of X is either 0 or 1.

- (a) (**5 points**) Explain what it means that an array is bitonic. Use at most three sentences.

Solution

Any of the solutions below (or similar) are acceptable:

- A bitonic sequence is either a monotonically increasing sequence followed by a decreasing sequence or the other way around.
- (Because I said X is all 0s and 1s,) A bitonic sequence is of the form $0^*1^*0^*$ or $1^*0^*1^*$.
- If a solution just says “increasing followed by decreasing” or “a sequence of 0s followed by a sequence of 1s followed by a sequence of 0s”, that will get full credit, even though the more general definition is needed for the proof in question 1c.

- (b) (**10 points**) Let Y be an array with N elements such that:

$$\begin{aligned} Y[i] &= \min(X[i], X[i + (N/2)]), & \text{if } 0 \leq i < N/2; \\ &= \max(X[i - (N/2)], X[i]), & \text{if } N/2 \leq i < N. \end{aligned}$$

Show that either

- $Y[i] = 0$ for $0 \leq i < N/2$, or
- $Y[i] = 1$ for $N/2 \leq i < N$.

Both may hold for particular choices of X , but you just need to show that at least one of these conditions holds.

Solution

- If X is all 0s, then Y is all 0s, and $Y[i] = 0$ for $0 \leq i < N/2$ which satisfies the claim.
- Otherwise, let i_{lo} and i_{hi} be defined as suggested in the hint for the problem. Note, that this assumes that X is increasing-then-decreasing bitonic. The argument if X is decreasing-then-increasing is equivalent. Consider the three cases below:
 - case $N/2 \leq i_{\text{lo}}$. In this case, $X[i] = 0$ for $0 \leq i < N/2 \leq i_{\text{lo}}$; $Y[i] = 0$ for $0 \leq i < N/2 \leq i_{\text{lo}}$; and the claim is satisfied.
 - case $i_{\text{hi}} < N/2$, this is analogous to the previous case. In particular, $X[i] = 1$ for $i_{\text{hi}} < N/2 \leq i < N \leq i_{\text{lo}}$; $Y[i] = 1$ for $i_{\text{hi}} < N/2 \leq i < N \leq i_{\text{lo}}$; and the claim is satisfied.
 - case $i_{\text{lo}} < N/2 \leq i_{\text{hi}}$. In this case, $Y[0..(N/2) - 1]$ is a sequence of 0s followed by a sequence of 1s, and $Y[(N/2)..(N - 1)]$ is a sequence of 1s followed by a sequence of 0s. This produces two cases depending on how their “breaks” line up:
 - * If $i_{\text{hi}} \leq i_{\text{lo}} + (N/2)$, then for all $0 \leq i < (N/2)$ either $X[i] = 0$ (i.e. for $0 \leq i < i_{\text{lo}}$) or $X[i + (N/2)] = 0$ (i.e. for $i_{\text{hi}} < i + (N/2) < N$). Therefore $Y[i] = 0$ for $0 \leq i < N/2$, and the claim is satisfied.
 - * Otherwise $i_{\text{lo}} + (N/2) < N$, and a similar argument shows that $Y[i] = 1$ for $N/2 \leq i < N$. In a bit more detail, for all $N/2 \leq i < N$, either either $X[i] = 1$ (i.e. for $N/2 \leq i \leq i_{\text{hi}}$) or $X[i - (N/2)] = 1$ (i.e. for $i_{\text{lo}} \leq i - (N/2) < N/2$). Therefore $Y[i] = 1$ for $N/2 \leq i < N$, and the claim is satisfied.

- (c) (**5 points**) Let Y be defined as in question 1b. Show that $Y[0..(N/2) - 1]$ is bitonic and that $Y[(N/2)..(N - 1)]$ is bitonic.

Solution

- If $(N/2) \leq i_{lo}$ then $X[0..(N/2)-1]$ is all 0s; $Y[0..(N/2)-1]$ is all 0s; and the claim holds.
- If $i_{hi} < (N/2)$ then $X[(N/2)..(N-1)]$ is all 0s; $Y[0..(N/2)-1]$ is all 0s; and the claim holds.
- The final case is $0 \leq i_{lo} < (N/2) \leq i_{hi} < N$. If $i_{lo} + (N/2) \leq i_{hi}$, then the 1s in the lower half of X overlap the 1s in the upper half of X , and $Y[0..(N/2)-1]$ is $0^*1^*0^*$ bitonic. Otherwise, $i_{lo} + (N/2) > i_{hi}$; the 0s in the lower half of X overlap the 0s in the upper half of X , $Y[0..(N/2)-1]$ is all 0s (and the top half of Y is $1^*0^*1^*$ bitonic).

- (d) **(5 points, Extra Credit)** Let A be a $N \times M$ array, where N is even, and $A[i, j]$ denotes the element of A in row i and column j . Let $A[i, 0..M-1]$ denote a row of A . Assume that every element of A is either 0 or 1, and that for $0 \leq i < N$, row $A[i, 0..M-1]$ is sorted into ascending order of i if i is even and into descending order if i is odd.

Let B be the $N \times M$ array where for $0 \leq j < M$, column $B[0..N-1, j]$ is $\text{sort}(A[0..N-1, j])$, where sort sorts the elements of the column into ascending order.

Prove that at least $N/2$ rows of B are either all 0s or all 1s. For example, if $N = 16$, and B has 6 rows that are all 0s and three rows that are all 1s, then the claim is satisfied.

Solution

Note that we could sort the columns of A by first doing a compare and swap of $A[2i, j]$ and $A[2i+1, j]$ for $0 \leq i < (N/2)$, and then sorting the result. Doing this for $0 \leq j < M$ is the same as doing a compare-and-swap of rows $2i$ and $2i+1$. From questions 1b and 1c this produces (at least) one row that is all 0s or all 1s. The sort will take that row to the bottom of B if the row is all 0s or to the top of B if it is all 1s. Thus, each pair of rows of A produces at least one row that is all 0s or all 1s in B . Because A has N rows, B must have at least $N/2$ rows that are all 0s or all 1s.

2. CUDA (20 points) Let's say that I have an array of

$$n = 2^{27} = 2^{17} \cdot 2^{10} = 131,072 \cdot 1,024 = 134,217,728$$

floats, and I want to compute their sum using my blazing-fast GPU.

- (a) **(10 points)** I've written three approaches to compute the sum below. Each is either incorrect or inefficient. Identify the problem with each, giving a short (at most three sentences) explanation for each one describing why it is wrong or why it is slow.

- i. Create one block with 134,217,728 threads executing the kernel:

```
0 __global__ void sum1(float *x, int n) {
1     // assume n is a power of 2
2     int myId = blockDim.x * blockIdx.x + threadIdx.x;
3     if(myId < n) {
4         for(int stride = 2; stride < n; stride *= 2) {
5             __syncthreads();
6             if((myId % stride) == 0)
7                 x[myId] = x[myId] + x[myId + stride]
8         }
9     }
10    // the result is in x[0]
11 }
```

The kernel launch is:

```
sum1<<<1, 134217728>>(x, 134217728)
```

- ii. The same kernel as for version [i], but this time the kernel launch is:

```
sum1<<<131072, 1024>>(x, 134217728)
```

- iii. Create one block with 1024 threads executing the kernel where each thread computes the sum of $n/1024$ elements before the threads perform a reduce:

```
0    __global__ void sum3(float *x, int n, int m) {
1        // assume n and m are powers of 2
2        int myId = blockDim.x * blockIdx.x + threadIdx.x;
3        if(myId < n) {
4            // compute the sum for my part of the array
5            float sum = 0.0;
6            for(int i = 0; i < m; i++)
7                sum += x[myId + i];
8
9            // reduce
10           for(int stride = 2; stride < blockDim.x; stride *= 2) {
11               __syncthreads();
12               if((myId % stride) == 0)
13                   x[myId] = x[myId] + x[myId + stride];
14           }
15       } // the result is in x[0]
16   }
```

The kernel launch is:

```
sum3<<<1, 1024>>(x, 134217728, 131072)
```

Solution

- Version (i) is incorrect because a block can have at most 1024 threads.
- Version (ii) is incorrect `__syncthreads()` only acts as a barrier for threads in the same block. It does not provide synchronization across blocks.
- Version (iii) is incorrect because
 - the threads compute their initial sums over overlapping parts of the array and miss most of the array.
 - the initial sum is not copied back to memory; so, it's not available for the reduce.

Note, these two errors were unintentional – oops. But, the problem lets you point out any error or efficiency issue; so, this just made the problem a little bit easier. No bug-bounty for these.

Version (iii) is also inefficient for many reasons (any one is fine):

- Thread divergence. In the reduce loop, half of the threads in a warp are active in the first iteration, one fourth in the second, and so on.
- Global memory accesses in the reduce loop. The global accesses in the initial sums are unavoidable, but we could use local memory here.
- Lots of synchronization: we could use “warp synchronous” execution once all of the active threads are part of the same warp.

- (b) **(10 points)** There are many ways to improve the performance of the slow kernel. Pick one such as coalescing global memory references, or using shared memory, or avoiding bank conflicts, or reducing thread divergence, or

- **(2 points)** State what optimization you are performing.
- **(4 points)** Write a short (no more than three sentence) description of what the optimization does.

- (2 points) Write the revised kernel. You can just write the lines that you change compared with the kernels I wrote above, and write

lines 05-08 of sum2

or similar to refer to lines from the kernels above.

Solution

I'll fix each of the problems listed above, but a real solution only needs to address one.

Fix sum3:

Change line 7 to:

```
sum += x[m*myId + i];
```

Just before the reduce (i.e. between lines 7 and 8) add

```
x[myId] = sum;
```

But now it will be slow due to uncoalesced memory references.

Global Memory Coalescing:

Change line 6 to:

```
int stride = blockDim.x * blockDim.x;
for(int i = myId; i < n; i += stride) {
```

Global memory accesses in the reduce:

Declare a shared array, for example between lines 1 and 2 add

```
__shared__ float y[1024];
```

Then, change the line I added between lines 7 and 8 to:

```
y[myId] = sum;
```

and change line 12 to

```
y[myId] = y[myId] + y[myId+stride];
```

For an even faster solution, each thread would continue to use `sum` until the last iteration before it becomes idle. It would copy `sum` to `y[myId]` in that last iteration.

Reducing thread divergence:

As described in lecture and the text, start stride at `blockDim.x/2` and divide it by two until it is 0. This changes line 9-12 to:

```
9      for(int stride = blockDim.x/2; stride > 0; stride /= 2) {
10         __syncthreads();
11         if((myId < stride) == 0)
12             y[myId] = y[myId] + y[myId + stride]
```

Using warp-synchronous execution:

Change the revised line 9 to:

```
for(int stride = blockDim.x/2; stride >= warpsize; stride /= 2) {
```

Then add another copy of the loop where `stride` goes from `warpsize/2` down to 0.

Omit the `__syncthreads()` in this second loop.

- (c) (10 points) Give two reasons that a GPU needs thousands of active threads to fully utilize the processors.
For each, you can give its name (e.g. the nVidia™ jargon word or phrase) and a short description (at most four or five sentences for each of your two reasons).

Solution

Here are more than two, but a real solution should just have two.

- **Mitigating data hazards:** Using lots of threads allows the SPs to have deep pipelines (easier to implement, lower power, higher clock speeds) without the complications of bypasses.

- **Mitigating control hazards:** Even though a branch instruction takes 10s of cycles (or more) on a GPU, the SM can fetch from other warps until the branch is resolved.
- **Hiding global memory latency:** Accesses to global memory are slow. Other warps can execute while a warp is waiting for a load to finish.
- **Warps have lots of threads:** This amortizes the hardware, and especially the power consumption for instruction fetch, decode, and other pipeline control issues across many SPs.
- **A GPU has lots of SMs:** This is how it gets lots of parallelism. But many SMs times many warps per SM times many threads per warp results in needing thousands of threads to get high performance from GPU.

3. GPUs (20 points)

Earlier this month, nVidia™ announced the GP100 GPU. It has 3584 single-precision floating point cores with a peak 10.6 teraflops of single-precision floating point computation (1 teraflop = 1000 gigaflops = 10^{12} floating point operations per second). The interface to off-chip, global memory has a bandwidth of 720 GBytes/sec.

- (a) (10 points) How many floating point operations must the GP100 perform on each `float` loaded from global memory if it is to achieve its peak floating point operation rate (and not be limited by memory bandwidth)?

Solution

A `float` is 4 bytes. Thus, the GP100 can load

$$= \frac{720 \frac{\text{Gbytes}}{\text{sec}} \cdot \frac{1\text{float}}{4\text{bytes}}}{180 \frac{\text{Gfloats}}{\text{sec}}}$$

As stated in the problem, The GP100 can perform 10.6 teraflops. The number of floating point operations we need to do per global memory access to allow for peak floating-point rate is:

$$= \frac{10.6 \cdot 10^{12} \frac{\text{float-ops}}{\text{sec}} \cdot \left(180 \cdot 10^9 \frac{\text{floats-from-gmem}}{\text{sec}}\right)^{-1}}{58.8 \frac{\text{float-ops}}{\text{float-from-gmem}}}$$

This is more than twice the CGMA needed to fully utilize the SPs on the GTX 550 Ti on the linXX boxes that we used.

- (b) (10 points) Give two reasons that a GPU needs thousands of active threads to fully utilize the processors. For each, you can give its name (e.g. the nVidia™ jargon word or phrase) and a short description (at most four or five sentences for each of your two reasons).

Solution

Here are more than two, but a real solution should just have two.

- **Mitigating data hazards:** Using lots of threads allows the SPs to have deep pipelines (easier to implement, lower power, higher clock speeds) without the complications of bypasses.
- **Mitigating control hazards:** Even though a branch instruction takes 10s of cycles (or more) on a GPU, the SM can fetch from other warps until the branch is resolved.
- **Hiding global memory latency:** Accesses to global memory are slow. Other warps can execute while a warp is waiting for a load to finish.
- **Warps have lots of threads:** This amortizes the hardware, and especially the power consumption for instruction fetch, decode, and other pipeline control issues across many SPs.
- **A GPU has lots of SMs:** This is how it gets lots of parallelism. But many SMs times many warps per SM times many threads per warp results in needing thousands of threads to get high performance from GPU.

4. **Speed-up** Let's say that I spend all day running four programs in sequence on my x86 linux machine:

- Program-A runs for 1 minute on the x86.
- Program-B runs for 1 minute on the x86.
- Program-C runs for 1 minute on the x86.
- Program-D runs for 1 minute on the x86.

Each program depends on the results of the one that came before it; so, I can't run them in parallel \odot . After Program-D finishes, I go back and start the sequence with Program-A again.

- (a) **(5 points)** Assuming that I work 8 hours per day and get paid \$1 each time I complete all four tasks, how much do I earn per day?

Solution

Completing all four tasks takes 4 minutes. That means 15 tasks in an hour, and 120 tasks in an eight-hour workday. I get paid \$120/day.

- (b) **(5 points)** An nVidiaTM sales person drops by my office and tells me that with a new, GP100 GPU, I can get the following speedups:

- Program-A: speed-up = 100.0;
- Program-B: speed-up = 10.0;
- Program-C: speed-up = 1.0;
- Program-D: speed-up = 0.1.

If I sell my linux machine and move all four of my computing tasks to the GPU, long will it take me to complete all four tasks? What is the overall speed-up?

Solution

The total time for all four tasks is

$$= \left(\frac{1}{100} + \frac{1}{10} + \frac{1}{1} + \frac{1}{0.1} \right) \text{ minutes}$$
$$= 11.11 \text{ minutes}$$

Running on the x86, the time for all four is 4 minutes. The speedup is

$$\begin{aligned} SpeedUp &= \frac{4\text{min}}{11.11\text{min}} \\ &= 0.36 \end{aligned}$$

Yikes! I've been swindled, and now I'll only get \$43/day. Ouch!

- (c) **(5 points)** Maybe selling my linux box wasn't such a good idea. Let's say I run Program-D on my linux machine, and I run the other three programs on the GPU. For simplicity, we'll ignore the time for communication when switching between running one task on the GPU and the next on the CPU (or the other way around). Because the tasks have sequential dependencies, I can't use the CPU and the GPU at the same time. With this combined CPU + GPU approach what is the speed-up compared with running everything on the CPU?

Solution

By calculation like the previous one, the total time is now for all four tasks is 2.11 minutes, and the speed-up is slightly less than 1.9.

On a happier note, I'll now earn \$211/day. Will the extra \$91/day cover the cost of the GP100? It should.

- (d) **(5 points)** What is Amdahl's law? Write the mathematical formula, and write two or three sentences to explain the intuition it gives for the speed-up results in the previous two sections.

Solution

The formula is:

$$\text{SpeedUp} = \frac{1}{f/s+(1-f)}$$

where $0 \leq f \leq 1$ is the fraction of an application that can be improved, and s is the speed-up for that fraction. In the limit that $s \rightarrow \infty$, the speed-up goes to $1/(1 - f)$. This is what we are seeing in the example above. Half of the tasks (programs C and D) don't benefit from the speed up. If the other two had really large speed-ups (and 100 isn't bad), then we would still only get a speed-up of 2. The speed-up of 1.896 is getting close to the value of 2. Note that if the "improvement" results in a slow-down for some portion (Program-D in the part a), then that can result in a very large performance degradation.

Simple version: it's tempting to just try to average the speed-ups. For the problem above the average is 27.775. That sounds great. In reality, we need to account for the slowest part. If we can adjust our workload to have most of it get the large speed-up, then we can get a huge pay-off. In practice, this means:

- Parallel computing doesn't solve all of your problems.
- We adjust our uses of computers to line up with the problems that do get large speed-ups, such as graphics, animation, simulations, data analytics, and machine learning. Applications that don't benefit from parallel computing may still be important, but they will evolve much more slowly.

5. Other stuff (20 points)

Answer each question below with 1-3 sentences. Points may be taken off for long answers.

- (a) **(4 points)** What is a pipeline bypass?

Solution

A mechanism that allows a result that is in the pipeline from an earlier issued instruction to be made available to a later issued instruction, even if the first instruction is still in the pipeline.

- (b) **(4 points)** What is a fused multiply-add?

Solution

A common feature in hardware for floating-point arithmetic where a multiply followed by an add, e.g. `a*x + y`, can be computed as single operation. This can be very useful for common numerical algorithms (e.g. matrix multiplication), and it lets the marketing people claim higher GFlop or TFlop numbers by counting the fused multiply-add as two floating point operations.

- (c) **(4 points)** What is the cross-section bandwidth of a message-passing multiprocessor?

Solution

Consider all ways to partition the N processors of the multiprocessor into two groups of size $N/2$. For each such partition, let B be the bandwidth of all the links between the two parts of the partition. The cross-section bandwidth is the smallest value of B for all possible partitions.

- (d) **(4 points)** What is a "straggler" in a map-reduce computation?

Solution

A straggler is a task that takes much longer than the others.

(e) **(4 points)** How do “back-up tasks” mitigate the problems caused by stragglers?

Solution

We can use “work stealing”. Near the end of a map-reduce (or other parallel) computation, assign incomplete tasks to idle processors. Often, the stragglers are due to processor being slow (other jobs running on it, configuration problems, etc.), failed, or network problems. By replicating the task, the computation can finish when either the original task or the back-up task completes.

Useful stuff

Amdahl's Law $speedup = \frac{1}{\frac{f}{s} + (1-f)}$

Erlang Cheat Sheet

lists:foldl

```
fsum(L) -> lists:foldl(fun(X,S) -> S+X end, 0, L).
```

reduce

```
largest_gap_seq([]) -> {};
largest_gap_seq([V]) -> {V, V};
largest_gap_seq([V1, V2 | Tail]) ->
    F = fun(XNew, {CurGap, XOld}) -> {gap_max(CurGap, {XOld, XNew}), XNew} end,
    {Gap, _XLast} = lists:foldl(F, {{V1, V2}, V2}, Tail),
    Gap.
```

```
largest_gap_par(WTree, DataKey) ->
    Leaf = fun(ProcState) -> leaf(workers:get(ProcState, DataKey)) end,
    Combine = fun(Left, Right) -> combine(Left, Right) end,
    Root = fun(Value) -> root(Value) end,
    wtree:reduce(WTree, Leaf, Combine, Root).
```

```
leaf([]) -> {};
leaf([V]) -> {V, {V,V}, V};
leaf(List) -> {hd(List), largest_gap_seq(List), lists:last(List)}.
```

```
combine({}, Right) -> Right;
combine(Left, {}) -> Left;
combine({L1, {L2, L3}, L4}, {R1, {R2, R3}, R4}) ->
    {L1, gap_max(gap_max({L2, L3}, {L4, R1}), {R2, R3}), R4}.
```

```
root({_MinV, {V1, V2}, _MaxV}) -> {V1, V2}.
```

```
gap_max({A1, A2}, {B1, B2}) when (A2 - A1) >= (B2 - B1) -> {A1, A2};
gap_max(_, B) -> B.
```

scan

```
bank_statement(W, SrcKey, DstKey, InitialBalance) ->
    Leaf1 = fun(ProcState) ->
        process_transactions(wtree:get(ProcState, SrcKey), {1, 0})
    end,
    Leaf2 = fun(ProcState, AccIn) ->
        {_S, C} = AccIn,
        Src = wtree:get(ProcState, SrcKey),
        Result = process_transactions_cum(Src, C), % compute the cumulative sum
        wtree:put(ProcState, DstKey, Result) % save the result -- must be the last expression
    end, % in the Leaf2 function
    Combine = fun({C1, S1}, {C2, S2}) -> {C1*C2, S1*C2 + S2} end,
    wtree:scan(W, Leaf1, Leaf2, Combine, {1, InitialBalance}).
```

cuda

```
--shared__ float v[1024];
__device__ float f(float x) {
    return ((5/2)*(x*x*x - x));
}

__device__ void compute_and_reduce(uint n, uint m, float *x) {
    uint myId = threadIdx.x;
    if(myId < n) {
        float y = x[myId];
        for(uint i = 0; i < m; i++) {
            y = f(y); % y = f^(i + 1) (x[myId])
        }
        v[myId] = y; % shared memory is much faster than global memory
        for(uint m = n >> 1; m > 0; m = n >> 1) { % reduce
            n -= m;
            __syncthreads();
            if(myId < m) {
                v[myId] += v[myId+n];
            }
        }
        x[myId] = v[myId]; % move result to global memory
    }
}
```

Parallel Computation

Mark Greenstreet

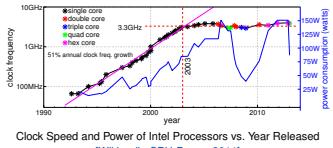
CpSc 418 – Jan. 4, 2017

Outline:

- Why Parallel Computation Matters
- Course Overview
- Our First Parallel Program
- The next month
- Table of Contents

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet
http://creativecommons.org/licenses/by/4.0/

Why Parallel Computation Matters



Clock Speed and Power of Intel Processors vs. Year Released [Wikipedia CPU-Power, 2011]

- In the good-old days, processor performance doubled roughly every 1.5 years.
- Single thread performance has seen small gains in the past 14 years.
- Too bad. If it had, we would have 1000GHz CPUs today. ☺
- Need other ways to increase performance.

Why Sequential Performance Can't Improve (much)

Power

- CPUs with faster clocks use more energy per operation than slower ones.
- For mobile devices: high power limits battery life.
- For desktop computers and gaming consoles: cooling high-power chips requires expensive hardware.
- For large servers and clouds, the power bill is a large part of the operating cost.

More Barriers to Sequential Performance

- The memory bottleneck.
 - ▶ Accessing main memory (i.e. DRAM) takes hundreds of clock cycles.
 - ▶ But, we can get high bandwidth.
- Limited instruction-level-parallelism.
 - ▶ CPUs already execute instructions in parallel.
 - ▶ But, the amount of this "free" parallelism is limited.
- Design complexity.
 - ▶ Designing a chip with 100 simple processors is way easier than designing a chip with one big processor.
- Reliability.
 - ▶ If a chip has 100 processors and one fails, there are still 99 good ones.
 - ▶ If a chip has 1 processor and it fails, then the chip is useless.
- See [Asanović et al., 2006].

Parallel Computers

- Mobile devices:
 - ▶ multi-core to get good performance on apps and reasonable battery life.
 - ▶ many dedicated "accelerators" for graphics, WiFi, networking, video, audio, ...
- Desktop computers
 - ▶ multi-core for performance
 - ▶ separate GPU for graphics
- Commercial servers
 - ▶ multiple, multi-core processors with shared memory.
 - ▶ large clusters of machines connected by dedicated networks.

Parallel Performance

The incentive for parallel computing is to do things that wouldn't be practical on a single processor.

- Performance matters.
- We need good models:
 - ▶ Counting operations can be very misleading – "adding is free."
 - ▶ Communication and coordination are often the dominant costs.
- We need to measure actual execution times of real programs.
 - ▶ There isn't a unified framework for parallel program performance analysis that works well in practice.
 - ▶ It's important to measure actual execution time and identify where the bottlenecks are.
- Key concepts with performance:
 - ▶ Amdahl's law, linear speed up, overheads.

Administrative Stuff – Who

- The instructors
 - ▶ **Mark Greenstreet**, mrg@cs.ubc.ca
 - ICCS 323, (604) 822-3065
 - Office hours: Tuesdays, 1pm – 2:30pm, ICCS 323
 - ▶ **Ian Mitchell**, mitchell@cs.ubc.ca
 - ICCS 217, (604) 822-2317
 - Office hours: Fridays, 12noon – 1pm, ICCS 217
- The TAs
 - ▶ **Devon Graham**, drgraham@cs.ubc.ca
 - ▶ **Chenxi Liu**, chenxiliu@cs.ubc.ca
 - ▶ **Carolyn Shen**, shen.carolyn@gmail.com
 - ▶ **Brenda Xiong**, krx.sky@gmail.com
- Course webpage: <http://www.ugrad.cs.ubc.ca/~cs418>.
- Online discussion group: on [piazza](#).

Homework

- Collaboration policy
 - ▶ You are welcome and encouraged to discuss the homework problems with other students in the class, with the TAs and me, and find relevant material in the text books, other book, on the web, etc.
 - ▶ You are expected to work out your own solutions and write your own code. Discussions as described above are to help understand the material. Your solutions must be your own.
 - ▶ You must properly cite your collaborators and any outside sources that you used. You don't need to cite material from class, the textbooks, or meeting with the TAs or instructor. See [slide 22](#) for more on the plagiarism policy.
- Late policy
 - ▶ Each assignment has an "early bird" date before the main date. Turn in your assignment by the early-bird date to get a 5% bonus.
 - ▶ **No late homework accepted.**

Grades: the big picture

$$\begin{aligned} \text{RawGrade} &= 0.35 \cdot HW + 0.25 \cdot MidTerm + 0.40 \cdot Final \\ \text{MiniBonus} &= 0.20 \cdot (1 - \min(\text{RawGrade}, 1)) \cdot \text{Mini} \\ BB &= 0.35 \cdot BB_{HW} + 0.25 \cdot BB_{MT} + 0.40 \cdot BB_{FX} \\ \text{CourseGrade} &= \min(\text{RawGrade} + \text{MiniBonus} + BB, 1) \times 100\% \end{aligned}$$

Lecture Outline

- Why Does Parallel Computation Matter?
- Course Overview
- Our First Parallel Program
 - ▶ Erlang quick start
 - ▶ Count 3s
 - ▶ Counting 3s in parallel
 - The root process
 - Spawning worker processes
 - The worker processes
 - Running the code

Outline

- Why Does Parallel Computation Matter?
- Course Overview
 - ▶ Topics
 - ▶ Syllabus
 - ▶ The instructor and TAs
 - ▶ The textbook(s)
 - ▶ Grades
 - Homework: 35% roughly one HW every two weeks
 - Midterm: 25% **March 1, in class**
 - Final: 40%
 - Mini-Assessments: see description on [slide 19](#)
 - Bug Bounties: see description on [slide 20](#)
 - ▶ Plagiarism – please don't
 - ▶ Learning Objectives
- Our First Parallel Program

Parallel Algorithms

- We'll explore some old friends in a parallel context:
 - ▶ Sum of the elements of an array
 - ▶ matrix multiplication
 - ▶ dynamic programming.
- And we'll explore some uniquely parallel algorithms:
 - ▶ Bitonic sort
 - ▶ mutual exclusion
 - ▶ producer consumer

Textbook(s)

- For Erlang: *Learn You Some Erlang For Great Good*, Fred Hebert,
 - ▶ Free! On-line at <http://learnyousomeerlang.com>.
 - ▶ You can buy the dead-tree edition at the same web-site if you like.
- For CUDA: *Programming Massively Parallel Processors: A Hands-on Approach* (2nd or 3rd ed.), D.B. Kirk and W.M.W. Hwu.
 - ▶ Please get a copy by late February – I'll assign readings starting after the midterm. It's available at amazon.ca and many other places.
- I'll hand-out copies of some book chapters:
 - ▶ *Principles of Parallel Programming* (chap. 5), C. Lin & L. Snyder – for the reduce and scan algorithms.
 - ▶ *An Introduction to Parallel Programming* (chap. 2), P.S. Pacheco – for a survey of parallel architectures.
 - ▶ Probably a few journal, magazine, or conference papers.

Exams

- Midterm, in class, on March 1.
- Final exam will be scheduled by the registrar.
- Both exams are open book, open notes, open homework and solutions – open anything printed on paper.
 - ▶ You can bring a calculator.
 - ▶ No communication devices: laptops, tablets, cell-phones, etc.

Plagiarism

- I have a very simple criterion for plagiarism:
 - ▶ Submitting the work of another person, whether be another student, something from a book, or something off the web and representing it as your own is plagiarism and constitutes academic misconduct.
- If the source is clearly cited, then it is not academic misconduct.
 - ▶ If you tell me "This is copied word for word from Jane Doe's solution" that is not academic misconduct. It will be graded as one solution for two people and each will get half credit. I guess that you could try telling me how much credit each of you should get, but I've never had anyone try this before.
- I encourage you to discuss the homework problems with each other.
 - ▶ If you're brainstorming with some friends and the key idea for a solution comes up, that's OK. In this case, add a note to your solution that lists who you collaborated with.
- More details at:
 - ▶ <http://www.ugrad.cs.ubc.ca/~cs418/plagiarism.html>
 - ▶ <http://www.ugrad.cs.ubc.ca/~cs418/2016-2/lecture/01-no-code.html>

Erlang Intro – very abbreviated!

- Erlang is a functional language:
 - ▶ Variables are given values when declared, and the value never changes.
 - ▶ The main data structures are lists, `[Head | Tail]`, and tuples (covered later).
 - ▶ Extensive use of pattern matching.
- The source code for the examples in this lecture is available at:
 - ▶ <http://www.ugrad.cs.ubc.ca/~cs418/2016-2/lecture/01-no-code.html>

Topics

- Parallel Architectures
- Parallel Performance
- Parallel Algorithms
- Parallel Programming Frameworks

Parallel Programming Frameworks

- Erlang: functional, message passing parallelism
 - ▶ Avoids many of the common parallel programming errors: races and side-effects.
 - ▶ You can write Erlang programs with such bugs, but it takes extra effort (esp. for the examples we consider).
 - ▶ Allows a simple presentation of many ideas.
 - ▶ But it's slow, for many applications, when compared with C or C++.
 - ▶ OTOH, it finds real use in large-scale distributed systems.
- CUDA: your graphics card is a super-computer
 - ▶ Excellent performance on the "right" kind of problem.
 - ▶ The data-parallel model is simple, and useful.

Parallel Barriers

Parallel Architectures

- **There isn't one, standard, parallel architecture for everything.**
 - ▶ We have:
 - ▶ Multi-core CPUs with a shared-memory programming model. Used for mobile device application processors, laptops, desktops, and many large data-base servers.
 - ▶ Networked clusters, typically running Linux. Used for web-servers and data-mining. Scientific supercomputers are typically huge clusters with dedicated, high-performance networks.
 - ▶ Domain specific processors
 - GPUs, video codecs, WiFi interfaces, image and sound processing, crypto engines, network packet filtering, and so on.
 - As a consequence, **there isn't one, standard, parallel programming paradigm.**

Why so many texts?

- There isn't one, dominant parallel architecture or programming paradigm.
- The Lin & Snyder book is a great, paradigm independent introduction.
- But, I've found that descriptions of real programming frameworks lack the details that help you write real code.
- So, I'm using several texts, but
 - ▶ You only have to buy one! ☺

Grades

Grades

- | | |
|-------------------|---|
| Homework: | 35% roughly one assignment every two weeks |
| Midterm: | 25% March 1, in class |
| Final: | 40% |
| Mini-Assessments: | see description on slide 19 |
| Bug Bounties: | see description on slide 20 |

Mini-Assessments

- Mini-assessments
 - ▶ Worth 20% of points missed from HW and exams.
 - ▶ If your raw grade is 90%, you can get at most 2% from the minis. Missing one or two isn't a big deal.
 - ▶ If your raw grade is 70%, you can get 6% from the minis. This can move your letter grade up a notch (e.g. C+ to B-).
 - ▶ If your raw grade is 45%, you can get up to 11% from the minis. Do the mini-assessments – I hate turning in failing grades.
 - ▶ The first is at <http://www.ugrad.cs.ubc.ca/~cs418/2016-2/min1/min1.pdf>, and due Jan. 9.
 - ▶ If you are on the course waitlist, we will select from the students who submit acceptable solutions to **Mini Assignment 1** to fill any slots that open.

Bug Bounties

- If I make a mistake when stating a homework problem, then the first person to report the error gets extra credit.
 - ▶ If the error would have prevented solving the problem, then the extra credit is the same as the value of the problem.
 - ▶ Smaller errors get extra credit in proportion to their severity.
- Likewise, bug bounties are awarded (as homework extra credit) for finding errors in mini-assessments, lecture slides, the course web-pages, code I provide, etc.
- The midterm and final have bug bounties awarded in midterm and final exam points respectively.
- If you find an error, report it.
 - ▶ Suspected errors in homework, lecture notes, and other course materials should be posted to piazza.
 - ▶ The first person to post a bug gets the bounty.
 - ▶ Bug-bounties reward you for looking at the HW when it first comes out, and not waiting until the day before it is due.

Learning Objectives (1/2)

- Parallel Algorithms
 - ▶ Familiar with parallel patterns such as reduce, scan, and tiling and can apply them to common parallel programming problems.
 - ▶ Can describe parallel algorithms for matrix operations, sorting, dynamic programming, and process coordination.
- Parallel Architectures
 - ▶ Can describe shared-memory, message-passing, and SIMD architectures.
 - ▶ Can describe a simple cache-coherence protocol.
 - ▶ Can identify how communication latency and bandwidth are limited by physical constraints in these architectures.
 - ▶ Can describe the difference between bandwidth and inverse latency, and how these impact parallel architectures.

Learning Objectives (2/2)

- Parallel Performance
 - ▶ Understands the concept of "speed-up": can calculate it from simple execution models or measured execution times.
 - ▶ Can identify key bottlenecks for parallel program performance including communication latency and bandwidth, synchronization overhead, and intrinsically sequential code.
- Parallel Programming Frameworks
 - ▶ Can implement simple parallel programs in Erlang and CUDA.
 - ▶ Can describe the differences between these paradigms.
 - ▶ Can identify when one of these paradigms is particularly well-suited (or badly suited) for a particular application.

Lists

- `[1, 4, 9, 16, 25, 36, 49, 64, 81, 100]` is a list of 10 elements.
- If `L` is a list, then `[0 | L]` is the list obtained by prepending the element `0` to the list `L`. In more detail:
 - ▶ `L = [1, 4, 9, 16, 25, 36, 49, 64, 81, 100]`
 - ▶ `2 >= [0 | L]`
 - ▶ `[0, 1, 4, 9, 16, 25, 36, 49, 64, 81, 100]`
 - ▶ `3 >= [0 | L]`
 - ▶ `[0, 1, 4, 9, 16, 25, 36, 49, 64, 81, 100]`
 - Of course, we traverse a list by using recursive functions:

Lists traversal example: sum

- ```
sum(List) ->
 if (length(List) == 0) -> 0;
 (length(List) > 0) -> hd(List) + sum(tl(List))
 end.
```
- Length (`L`) returns the number of elements in list `L`.
  - Head (`hd(L)`) returns the first element of list `L` (the head), and throws an exception if `L` is the empty list.
  - Tail (`tl(L)`) returns the list of all elements after the first (the tail).
  - `sum([1, 2, 3]) = [2, 3]. tl([1]) = []`
  - See `sum_wo_pm` ("sum without pattern matching") in [simple.erl](#)

## Pattern Matching – first example

We can use Erlang's pattern matching instead of the `if` expression:

```
sum([]) -> 0;
sum([Head | Tail]) -> Head + sum(Tail).
```

`sum([Head | Tail])` matches any non-empty list with `Head` being bound to the value of the first element of the list, and `Tail` being bound to the list of all the other elements.

More generally, we can use patterns to identify the different cases for a function.

This can lead to very simple code where function definitions follow the structure of their arguments.

See [sum in simple.erl](#)

## Count 3's: a simple example

Given an array (or list) with  $N$  items, return the number of those elements that have the value 3.

```
count3s([]) -> 0;
count3s([_| Tail]) -> 1 + count3s(Tail);
count3s([_ | Tail]) -> count3s(Tail).
```

We'll need to put the code in an erlang module. See `count3s` in [count3s.erl](#) for the details.

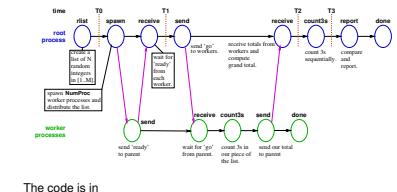
To generate a list of random integers, `count3s.erl` uses the function `rlist(N, M)` from course Erlang library that returns a list of  $N$  integers randomly chosen from  $1..M$ .

## Running Erlang

```
bash-3.2$ erl
Erlang/OTP 18 [erts-7.0] [source] ...
Eshell V7.0 (abort with ^G)

1> c(count3s).
{ok, count3s}
2> L20 = count3s:rlist(20,5).
[3, 4, 5, 3, 2, 3, 5, 4, 3, 3, 1, 2, 4, 1, 3, 2, 3, 3, 1, 3]
3> count3s:count3s(L20).
9
4> count3s:count3s(count3s:rlist(1000000,10)).
99961
5> q().
o
6> bash-3.2$
```

## A Parallel Version



The code is in

<http://www.ugrad.cs.ubc.ca/~cs418/2016-2/lecture/01-04/src/count3s.erl>

## Preview of the next month

January 6: Introduction to Erlang Programming  
Reading: [Learn You Some Erlang](#), the first eight sections – [Introduction through Recursion](#). Feel free to skip the stuff on [bit syntax](#) and [higher-order functions](#).

January 8: Processes and Messages  
Reading: [Learn You Some Erlang](#), Higher Order Functions and [The Hitchhiker's Guide...](#), through More on Multiprocessing  
Homework: [Homework 1 deadline out \(due Jan. 18\)](#) – Erlang programming  
Mini-Assignment: [Mini-Assignment 1 – 10:00am](#)

January 11: Reduce  
Reading: [Learn You Some Erlang](#), Errors and Exceptions through [A Short Visit to Common Data Structures](#)

January 13: Scan  
Reading: Lin & Snyder, chapter 5, pp. 112–125  
Homework: [Homework 2 deadline out \(due Jan. 13\)](#)

January 16: Generalized Reduce and Scan  
Homework: Homework 1 deadline for early-bird bonus (11:59pm)  
[Homework 2 goes out \(due Feb. 1\)](#) – Reduce and Scan

January 18: Reduce and Scan Examples  
Homework: Homework 2 deadline (due 11:59pm)

January 20-27: Parallel Architecture

January 28-February 6: Parallel Performance

## Review Questions

- Name one, or a few, key reasons that parallel programming is moving into mainstream applications.
- How does the impact of your mini assignment total on your final grade depend on how you did on the other parts of the class?
- What are bug-bounds?
- What is the count's problem?
- How did we measure running times to compute speed up?
  - Why did one approach show a speed-up greater than the number of cores used?
  - Why did the other approach show that the parallel version was **slower** than the sequential one?

## Supplementary Material

- [Erlang Resources](#)
- [Bibliography](#)
- [Table of Contents](#) – at the end!!!

## Erlang Resources

- Learn You Some Erlang  
<http://learnyousomeerlang.com>  
An on-line book that gives a very good introduction to Erlang. It has great answers to the “Why is Erlang this way?” kinds of questions, and it gives realistic assessments of both the strengths and limitations of Erlang.
- Erlang Examples:  
<http://www.ugrad.cs.ubc.ca/~cs418/2012-1/lecture/09-08.pdf>  
My lecture notes that walk through the main features of Erlang with examples for each. Try it with an Erlang interpreter running in another window so you can try the examples and make up your own as you go. This will cover everything you’ll need to make it through all (or most) of what we’ll do in class, but it doesn’t explain how to think in Erlang as well as “Learn You Some Erlang” or Armstrong’s Erlang book (next slide).

## More Erlang Resources

- The erlang.org tutorial!  
[http://www.erlang.org/doc/getting\\_started/users\\_guide.html](http://www.erlang.org/doc/getting_started/users_guide.html)  
Somewhere between my “Erlang Examples” and “Learn You Some Erlang”
- Erlang Language Manual  
[http://www.erlang.org/doc/reference\\_manual/users\\_guide.html](http://www.erlang.org/doc/reference_manual/users_guide.html)  
My go-to place when looking up details of Erlang operators, etc.
- On-line API documentation:  
<http://www.erlang.org/erldoc>
- The book: *Programming Erlang: Software for a Concurrent World*, Joe Armstrong, 2007,  
<http://pragprog.com/book/jaerlang/programming-erlang>  
Very well written, with lots of great examples. More than you'll need for this class, but great if you find yourself using Erlang for a big project.
- More resources listed at <http://www.erlang.org/doc.html>.

## Getting Erlang

- You can run Erlang by giving the command `erl` on any departmental machine. For example:
  - Linux: bowen, thetis, lin01, ..., lin25, ..., all machines above are `.ugrad.cs.ubc.ca`, e.g. `bowen.ugrad.cs.ubc.ca`, etc.
- You can install Erlang on your computer
  - Erlang solutions provides packages for Windows, OSX, and the most common Linux distros  
<http://erlang-solutions.com/resources/download.html>
  - Note: some linux distros come with Erlang pre-installed, but it might be an old version. You should probably install from the link above.

## Starting Erlang

- Start the Erlang interpreter.  
`thisis% erl`  
Erlang/OTP 18 [erts-7.0] [source] ...  
Eshell V7.0 (abort with ^G)
- 1> 2+3.  
5  
2>
- The Erlang interpreter evaluates expressions that you type.
- Expressions end with a “.” (period).

## Biography

- Krste Asanovic, Ras Bodik, et al.  
The landscape of parallel computing research: A view from Berkeley. Technical Report UCB/ECECS-2006-183, Electrical Engineering and Computer Science Department, University of California, Berkeley, December 2006.  
<http://www.eecs.berkeley.edu/Pubs/TechRpts/2006/ECCS-2006-183.pdf>
- Microprocessor quick reference guide.  
<http://www.intel.com/pressroom/kits/quickrefry.htm>, June 2013.  
accessed 29 August 2013.
- List of CPU power dissipation.  
[http://en.wikipedia.org/wiki/List\\_of\\_CPU\\_power\\_dissipation](http://en.wikipedia.org/wiki/List_of_CPU_power_dissipation), April 2011.  
accessed 26 July 2011.

## Table of Contents (1/2)

- Motivation
- Course Overview
  - Topics
    - Computer Architecture
    - Performance Analysis
    - Algorithms
    - Languages, Paradigms, and Frameworks
  - Syllabus
  - Course Administration – who's who
  - The Textbook(s)
  - Grades
    - Homework
    - Midterm and Final Exams
    - Mini-Assessments
    - Bug Bounties
  - Plagiarism Policy
  - Learning Objectives

## Table of Contents (2/2)

- Our First Parallel Program
  - Introduction to Erlang
  - The Count 3s Example
- Preview of the next month
- Review of this lecture
- Supplementary Material
  - Erlang Resources
  - Bibliography
  - Table of Contents

## Introduction to Erlang

Mark Greenstreet  
CpSc 418 – January 6, 2016

## Objectives

- Learn/review key concepts of functional programming:
  - Referential transparency.
  - Structuring code with functions.
- Introduction to Erlang
  - Basic data types and operations.
  - Program design by structural decomposition.
  - Writing and compiling an Erlang module.

## Erlang Basics

- Numbers:
  - Numerical Constants: 1, #8311, 1.5, 1.5e3,
  - but not: 1, .0, 5
- Arithmetic: +, -, \*, /, div, band, bor, bnot, bal, bxor
- Booleans:
  - Comparisons: ==, /=, ==/=, /=, <, <=, >, >=
  - Boolean operations (strict): and, or, not, xor
  - Boolean operations (short-circuit): andalso, orelse
- Atoms:
  - Constants: x, 'big DOG'
  - Operations: tests for equality and inequality. Therefore pattern matching.

## Lists and Tuples

- Lists:
  - Construction: [1, 2, 3], [Element1, Element2, ..., ElementN | Tail]
  - Operations: hd, tl, length, ++, -
  - Erlang's list library: <http://erlang.org/doc/man/lists.html>: all, any, filter, foldl, foldr, map, nth, nthtail, seq, sort, split, zipwith, and many more.
- tuples:
  - Construction: {1, dog, "called Rover"}
  - Operations: element, setelement, tuple\_size.
  - Lists vs. Tuples:
    - Lists are typically used for an **arbitrary** number of elements of the same “type” – like arrays in C, Java, ...
    - Tuples are typically used for a **fixed** number of elements of the varying “types” – like a struct in C or an object in Java.

## Strings

Mark Greenstreet  
CpSc 418 – January 6, 2016

## Erlang Makes Parallel Programming Easier

- Erlang is functional
  - Each variable gets its value when it's declared – it **never** changes.
  - Erlang eliminates many kinds of races – another process **can't** change the value of a variable while you're using it, because the values of variables never change.
- Erlang uses message passing
  - Interactions between processes are under explicit control of the programmer.
  - Fewer races, synchronization errors, etc.
- Erlang has simple mechanisms for process creation and communication
  - The structure of the program is not buried in a large number of calls to a complicated API.

Big picture: Erlang makes the issues of parallelism in parallel programs more apparent and makes it easier to avoid many common pitfalls in parallel programming.

## Example: Sorting a List

- The simple cases:
  - Sorting an empty list: `sort([]) -> []`
  - Sorting a singleton list: `sort([A]) -> [A]`
- How about a list with more than two elements?
  - Merge sort?
  - Quick sort?
  - Bubble sort? (NO WAY! Bubble sort is DISGUSTING!!)
- Let's figure it out.

## Referential Transparency

- This notion that a variable gets a value when it is declared and that the value of the variable never changes is called **referential transparency**.
  - In other words, the term many times in class – I thought it would be a good idea to let you know what it means. ☺
- We say that the value of the variable is **bound** to the variable.
- Variables in functional programming are much like those in mathematical formulas:
  - If a variable appears multiple places in a mathematical formula, we assume that it has the same value everywhere.
  - This is the case in a functional program.
  - This is **not** the case in an imperative program. We can declare `x` on line 17; assign it a value on line 20; and assign it another value on line 42.
  - The value of `x` when executing line 21 is different than when executing line 43.

## Loops violate referential transparency

```
// vector dot-product
sum = 0;
for(i = 0; i < a.length; i++) {
 if(a[i] == b[i]) {
 last = a[i];
 a = a.next;
 last->next = null;
 } else {
 ...
 }
}
```

- Loops rely on changing the values of variables.
- Functional programs use recursion instead.
- See also the [LYSE explanation](#).

## Life without loops

- Use recursive functions instead of loops:
  - `dotProd([], []) -> 0;`
  - `dotProd([A | At], [B | Bt]) -> A*B + dotProd(At, Bt).`
- Functional programming uses recursion instead of iteration:
  - `dotProd([], []) -> 0;`
  - `dotProd([A | At], [B | Bt]) -> A*B + dotProd(At, Bt).`
- Anything you can do with iteration can be done with recursion.
  - But the converse is not true (without dynamically allocating data structures).
  - Example: tree traversal.

## split (L)

Identify the cases and their return values according to the shape of `L`:

```
% If L is empty (recall that split returns a tuple of two lists):
split([]) -> { , }

% If L
split() -> { , }

% If L
split() -> { , }
```

## Finishing merge sort

- An exercise for the reader – see [slide 29](#).
- Sketch:
  - Write `merge(List1, List2) -> List12` – see [slide 30](#)
  - Write an Erlang module with the `sort`, `split`, and `merge` functions – see [slide 31](#)
  - Run the code – see [slide 33](#)

|                                                                                                                                                                                                                                                                                                                                                                                                                                                              |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------|---------|
| Fun with functions                                                                                                                                                                                                                                                                                                                                                                                                                                           |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Programming with Patterns                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |                                                                                                                                           | Anonymous Functions                                                                                                                                                     |                                                                                                                            | Higher-Order Functions                                                                                                            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Programming with patterns <ul style="list-style-type: none"><li>▷ often the code just matches the shape of the data</li><li>▷ like CPSC 110, but pattern matching makes it obvious</li><li>▷ see <a href="#">slide 16</a></li></ul>                                                                                                                                                                                                                        | • Fun expressions <ul style="list-style-type: none"><li>▷ in-line function definitions</li><li>▷ see <a href="#">slide 17</a></li></ul>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      | • Higher-order functions <ul style="list-style-type: none"><li>▷ encode common control-flow patterns</li><li>▷ see <a href="#">slide 18</a></li></ul>                                                                                                                                                                                                                                                                                                                                                                                                                                     | • List comprehensions <ul style="list-style-type: none"><li>▷ common operations on lists</li><li>▷ see <a href="#">slide 19</a></li></ul> | • Tail call elimination <ul style="list-style-type: none"><li>▷ makes recursion as fast as iteration (in simple cases)</li><li>▷ see <a href="#">slide 20</a></li></ul> | • Let's try it <ul style="list-style-type: none"><li>▷ examples:leafCount([1, 2, [3, 4, []], [5, [6, banana]]]).</li></ul> | • Notice how we used <b>patterns</b> to show how the recursive structure of <code>leafCount</code> follows the shape of the tree. | • Factorial = % We can even write recursive fun expressions! <pre>fun(X, Y) -&gt; X*X + Y*Y end.<br/>#Fun{level=12,52032458} =&gt; ok, but what can I do with it??<br/>4&gt; F = fun(X, Y) -&gt; X*X + Y*Y end.<br/>#Fun{level=12,52032458}<br/>5&gt; F(3, 4).<br/>25<br/>6&gt; Factorial = % We can even write recursive fun expressions!<br/>fun(F) -&gt; 1;<br/>F(N) when is_integer(N), N &gt; 0 -&gt; N*Fact(N-1)<br/>end.<br/>7&gt; Factorial(3).<br/>9<br/>8&gt; Factorial(3).<br/>* 1: variable 'Fact' is unbound<br/>9&gt; Factorial(-2).<br/>** exception error: no function clause matching<br/>erlang!'-inside-an-interpreted-fun' (-2)<br/>10&gt; Factorial(frog).<br/>** exception error: function clause matching<br/>erlang!'-inside-an-interpreted-fun-' (frog)<br/>11&gt;</pre> |                       |         |
| Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                             | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       | CS 418 – Jan. 6, 2016                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | 15 / 26                                                                                                                                   | Mark Greenstreet                                                                                                                                                        | Introduction to Erlang                                                                                                     | CS 418 – Jan. 6, 2016                                                                                                             | 16 / 26                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                       |         |
| List Comprehensions                                                                                                                                                                                                                                                                                                                                                                                                                                          |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | Introduction to Erlang                                                                                                                    | CS 418 – Jan. 6, 2016                                                                                                                                                   | 16 / 26                                                                                                                    | Mark Greenstreet                                                                                                                  | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            | CS 418 – Jan. 6, 2016 | 17 / 26 |
| Head vs. Tail Recursion                                                                                                                                                                                                                                                                                                                                                                                                                                      |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Anonymous Functions                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       |                                                                                                                                           | Higher-Order Functions                                                                                                                                                  |                                                                                                                            | Tail Call Elimination – a few more notes                                                                                          |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Map and filter are such common operations, that Erlang has a simple syntax for such operations.                                                                                                                                                                                                                                                                                                                                                            | • Both grow linearly for $n \leq 10^6$ <ul style="list-style-type: none"><li>▷ The tail recursive version has runtimes about 2/3 of the head-recursive version.</li></ul>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • It's called a <b>List Comprehension</b> : <ul style="list-style-type: none"><li>▷ <code>[Expr    Var &lt;- List, Cond, ...]</code>.</li><li>▷ <code>Expr</code> is evaluated with <code>Var</code> set to each element of <code>List</code> that satisfies <code>Cond</code>.</li><li>▷ Example:<pre>13&gt; L = [count3s:rlist(S, 1000).<br/>[446, 724, 946, 502, 312].<br/>14&gt; L[X*X    X &lt;- R, X rem 3 == 0].<br/>[197136, 97344].</pre></li></ul> | • For $N > 10^6$ , <ul style="list-style-type: none"><li>▷ The tail recursive version continues to have run-time linear in <math>N</math>.</li><li>▷ The head recursive version becomes much slower than the tail recursive version.</li></ul>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • See also <a href="#">List Comprehensions</a> in <a href="#">LYSE</a> .                                                                                                                                                                                                                                                                                                                                                                                     | • The Erlang compiler optimizes tail calls <ul style="list-style-type: none"><li>▷ When the last operation of a function is to call another function, the compiler just revises the current stack frame and jumps to the entry point of the callee.</li><li>▷ The compiler has turned the recursive function into a while-loop.</li><li>▷ Conclusion: When people tell you that recursion is slower than iteration – don't believe them!</li></ul>                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                             | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       | CS 418 – Jan. 6, 2016                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | 19 / 26                                                                                                                                   | Mark Greenstreet                                                                                                                                                        | Introduction to Erlang                                                                                                     | CS 418 – Jan. 6, 2016                                                                                                             | 20 / 26                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                       |         |
| Summary                                                                                                                                                                                                                                                                                                                                                                                                                                                      |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | Introduction to Erlang                                                                                                                    | CS 418 – Jan. 6, 2016                                                                                                                                                   | 20 / 26                                                                                                                    | Mark Greenstreet                                                                                                                  | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            | CS 418 – Jan. 6, 2016 | 21 / 26 |
| Review Questions                                                                                                                                                                                                                                                                                                                                                                                                                                             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | A Few More Review Questions                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |                                                                                                                                           | Review Questions                                                                                                                                                        |                                                                                                                            | Review Questions                                                                                                                  |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Why Erlang? <ul style="list-style-type: none"><li>▷ Functional – avoid complications of side-effects when dealing with concurrency.</li><li>▷ But, we can't use imperative control flow constructions (e.g. loops).</li><li>▷ Design by declaration: look at the structure of the data.</li><li>▷ More techniques coming in upcoming lectures.</li></ul>                                                                                                   | • What is the difference between <code>==</code> and <code>=:=</code> ?                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Sequential Erlang <ul style="list-style-type: none"><li>▷ Lists, tuples, atoms, expressions.</li><li>▷ Using structural design to write functions: example sorting.</li><li>▷ Functions: patterns, higher-order functions, head vs. tail recursion.</li></ul>                                                                                                                                                                                              | • Which of the following are valid Erlang variables, atoms, both, or neither? <pre>Foo, foo, 25, '25', 'Foo foo',<br/>'4 score and 7 years ago', X2,<br/>'4 score and 7 years ago'.</pre>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                             | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       | CS 418 – Jan. 6, 2016                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | 23 / 26                                                                                                                                   | Mark Greenstreet                                                                                                                                                        | Introduction to Erlang                                                                                                     | CS 418 – Jan. 6, 2016                                                                                                             | 24 / 26                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                       |         |
| Supplementary Material                                                                                                                                                                                                                                                                                                                                                                                                                                       |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | Introduction to Erlang                                                                                                                    | CS 418 – Jan. 6, 2016                                                                                                                                                   | 24 / 26                                                                                                                    | Mark Greenstreet                                                                                                                  | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            | CS 418 – Jan. 6, 2016 | 25 / 26 |
| Erlang Resources                                                                                                                                                                                                                                                                                                                                                                                                                                             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Finishing the merge sort example                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          |                                                                                                                                           | merge (L1, L2)                                                                                                                                                          |                                                                                                                            | Remarks about Constructing Lists                                                                                                  |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| The remaining material is included in the web-version of these slides:<br><a href="http://www.ugrad.cs.ubc.ca/~cs418/2016-2/lecture/01-06/slides.pdf">http://www.ugrad.cs.ubc.ca/~cs418/2016-2/lecture/01-06/slides.pdf</a>                                                                                                                                                                                                                                  | • Install Erlang on your computer <ul style="list-style-type: none"><li>▷ Erlang solutions provides packages for Windows, OSX, and the most common linux distros</li><li>▷ Note: some linux distros come with Erlang pre-installed, but it might be an old version. You should probably install from the link above.</li></ul>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               | • Write <code>merge(List1, List2) -&gt; List12</code> – see <a href="#">slide 30</a>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| I'm omitting it from the printed handout to save a few trees.                                                                                                                                                                                                                                                                                                                                                                                                | • The module declaration is a line of the form: <pre>-module(modulename).</pre>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | • Write an Erlang module with the <code>sort</code> , <code>split</code> , and <code>merge</code> functions – see <a href="#">slide 31</a>                                                                                                                                                                                                                                                                                                                                                                                                                                                |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Erlang resources.                                                                                                                                                                                                                                                                                                                                                                                                                                          | • Function exports are written as: <pre>-export([functionName/arity1,<br/>functionName/arity2, ...]).</pre>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  | • Run the code – see <a href="#">slide 33</a>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Finishing the merge sort example.                                                                                                                                                                                                                                                                                                                                                                                                                          | • The list of functions may span multiple lines and there may be more than one <code>-export</code> attribute: <pre>arity is the number of arguments that the function has. For example, if we define<br/>foo(A, B) -&gt; A+B.<br/>Then we could export foo with<br/>-export([..., foo/2, ...]).</pre>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       | Let's try it!                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Common mistakes with lists and how to avoid them.                                                                                                                                                                                                                                                                                                                                                                                                          | • There are many other attributes that a module can have. We'll skip the details. If you really want to know, it's all described <a href="#">here</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      | <pre>1&gt; c(sort).<br/>ok<br/>2&gt; R20 = count3s:rlist(20, 100). % test case: a random list<br/>[45, 73, 95, 51, 32, 60, 67, 48, 60, 15, 21, 70, 16, 56, 22, 46, 43, 1, 57]<br/>3&gt; S20 = sort(r20). % sort it<br/>[1, 15, 16, 21, 22, 32, 43, 45, 46, 48, 51, 56, 57, 60, 60, 67, 70, 73, 92, 95]<br/>4&gt; R20 -- S20. % empty if each element in R20 is in S20<br/>[]<br/>5&gt; S20 -- R20. % empty if each element in S20 is in R20<br/>[]</pre>                                                                                                                                  |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • A few remarks about atoms.                                                                                                                                                                                                                                                                                                                                                                                                                                 | • Rules of Erlang punctuation: <ul style="list-style-type: none"><li>▷ Erlang declarations end with a period: .</li><li>▷ A declaration can consist of several alternatives.</li><li>▷ Alternatives are separated by a semicolon ;</li><li>▷ Note that many Erlang constructions such as <code>case</code>, <code>fun</code>, <code>if</code>, and <code>receive</code> have multiple alternatives as well.</li><li>▷ A declaration or alternative can be a block expression<ul style="list-style-type: none"><li>▷ Expressions in a block are separated by a comma ,</li><li>▷ The value of a block expression is the last expression of the block.</li></ul></li><li>▷ Expressions that begin with a keyword end with end<ul style="list-style-type: none"><li>▷ case Alternatives end</li><li>▷ Fun Alternatives end</li><li>▷ If Alternatives end</li><li>▷ receive Alternatives end</li></ul></li></ul> | • Yay – it works!!! (for one test case)                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Suppressing verbose output when using the Erlang shell.                                                                                                                                                                                                                                                                                                                                                                                                    | • Erlang rules of punctuation:                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               | • The code is available at <a href="http://www.ugrad.cs.ubc.ca/~cs418/2016-2/lecture/01-06/src/sort.erl">http://www.ugrad.cs.ubc.ca/~cs418/2016-2/lecture/01-06/src/sort.erl</a>                                                                                                                                                                                                                                                                                                                                                                                                          |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Forgetting variable bindings (only in the Erlang shell).                                                                                                                                                                                                                                                                                                                                                                                                   | • A few other things to note about punctuation:                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Avoiding Verbose Output                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Table of Contents                                                                                                                                                                                                                                                                                                                                                                                                                                          | • Atom constants are special constants.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      | Sometimes, when using Erlang interactively, we want to declare a variable where Erlang would spew enormous amounts of "uninteresting" output were it to print the variable's value.                                                                                                                                                                                                                                                                                                                                                                                                       |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                             | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       | CS 418 – Jan. 6, 2016                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | 31 / 26                                                                                                                                   | Mark Greenstreet                                                                                                                                                        | Introduction to Erlang                                                                                                     | CS 418 – Jan. 6, 2016                                                                                                             | 32 / 26                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                       |         |
| Misconstructing Lists                                                                                                                                                                                                                                                                                                                                                                                                                                        |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | Introduction to Erlang                                                                                                                    | CS 418 – Jan. 6, 2016                                                                                                                                                   | 32 / 26                                                                                                                    | Mark Greenstreet                                                                                                                  | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            | CS 418 – Jan. 6, 2016 | 33 / 26 |
| Punctuation                                                                                                                                                                                                                                                                                                                                                                                                                                                  |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Remarks about Atoms                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       |                                                                                                                                           | Objectives                                                                                                                                                              |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Working with <code>divisible_drop</code> from the previous slide...                                                                                                                                                                                                                                                                                                                                                                                          | • Erlang has lots of punctuation: commas, semicolons, periods, and end.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      | • Atom is a special constant.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | • Introduce Erlang's features for concurrency and parallelism                                                                             |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Now, change the second alternative in the if to<br><pre>A rem N == 0 -&gt; [A, divisible_filter(N, Tail)]</pre>                                                                                                                                                                                                                                                                                                                                            | • It's easy to get syntax errors or non-working code by using the wrong punctuation somewhere.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               | • Atoms can be compared for equality.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | ▷ Spawning processes.                                                                                                                     |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Try the previous test case:<br><pre>?&gt; examples:divisible_drop(3, [0, 1, 4, 17, 42, 100]).<br/>[1, 4, 17, [100, []]]</pre>                                                                                                                                                                                                                                                                                                                              | • Rules of Erlang punctuation: <ul style="list-style-type: none"><li>▷ A declaration can consist of several alternatives.</li><li>▷ Alternatives are separated by a semicolon ;</li><li>▷ Note that many Erlang constructions such as <code>case</code>, <code>fun</code>, <code>if</code>, and <code>receive</code> have multiple alternatives as well.</li><li>▷ A declaration or alternative can be a block expression<ul style="list-style-type: none"><li>▷ Expressions in a block are separated by a comma ,</li><li>▷ The value of a block expression is the last expression of the block.</li></ul></li><li>▷ Expressions that begin with a keyword end with end<ul style="list-style-type: none"><li>▷ case Alternatives end</li><li>▷ Fun Alternatives end</li><li>▷ If Alternatives end</li><li>▷ receive Alternatives end</li></ul></li></ul>                                                    | • Actually, any two Erlang can be compared for equality, and any two terms are ordered.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   | ▷ Sending and receiving messages.                                                                                                         |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Moral: If you see a list that is nesting way too much, check to see if you wrote a comma where you should have used a ;.                                                                                                                                                                                                                                                                                                                                   | • Erlang declarations end with a period: .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   | • Each atom is unique.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    | ▷ Describe timing measurements for these operations and the implications for writing efficient parallel programs.                         |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Restore the code and then change the second alternative for <code>divisible_drop</code> to <code>divisible_drop(N, [A, Tail])</code><br>→ Try running our previous test:<br><pre>?&gt; examples:divisible_drop(3, [0, 1, 4, 17, 42, 100]).<br/>** exception error: no function clause matching...</pre>                                                                                                                                                    | • A declaration can consist of several alternatives.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         | • Syntax of atoms <ul style="list-style-type: none"><li>▷ Anything that looks like an identifier and starts with a lower-case letter, e.g. x.</li><li>▷ Anything that is enclosed between a pair of single quotes, e.g. '4 BIG apples'.</li><li>▷ Some languages (e.g. Matlab or Python) use single quotes to enclose string constants, some (e.g. C or Java) use double quotes to enclose character constants.</li><li>▷ It's not an Erlang.</li><li>▷ The atom '47 big apples' is not a string or a list, or a character constant.</li><li>▷ It's just its own, unique value.</li></ul> | ▷ Communication often dominates the runtime of parallel programs.                                                                         |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • In the Erlang shell, f( <code>Variable</code> ) makes the shell "forget" the binding for the variable.                                                                                                                                                                                                                                                                                                                                                     | • Atom constants are written with single quotes, but they are not strings.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   | • Atom constants are special constants.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   | • The source code for the examples in this lecture is available here: <a href="#">procs.erl</a> .                                         |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| <pre>1&gt; X = 2+3.<br/>5.<br/>13&gt; X = 2+3.<br/>** exception error: no match of right hand side value 6.<br/>14&gt; f(X).<br/>ok<br/>15&gt; X = 2+3.<br/>6<br/>16&gt;</pre>                                                                                                                                                                                                                                                                               | • Erlang Basics – basic types and their operations.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | • Processes                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Forgetting Bindings                                                                                                                                                                                                                                                                                                                                                                                                                                          | • Fun with functions – patterns, anonymous functions, higher-order functions, list comprehensions, head vs. tail recursion                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   | • Messages                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Referential transparency means that bindings are forever.                                                                                                                                                                                                                                                                                                                                                                                                  | • Preview of upcoming lectures                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               | • Timing Measurements                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • This can be nuisance when using the Erlang shell.                                                                                                                                                                                                                                                                                                                                                                                                          | • Review of this lecture                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | • Preview, Review, etc.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • Sometimes we assign a value to a variable for debugging purposes.                                                                                                                                                                                                                                                                                                                                                                                          | • Supplementary Material                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | • Table of Contents                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • We'd like to overwrite that value later so we don't have to keep coming up with more name.s                                                                                                                                                                                                                                                                                                                                                                |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| • In the Erlang shell, f( <code>Variable</code> ) makes the shell "forget" the binding for the variable.                                                                                                                                                                                                                                                                                                                                                     |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| <pre>1&gt; X = 2+3.<br/>5.<br/>13&gt; X = 2+3.<br/>** exception error: no match of right hand side value 6.<br/>14&gt; f(X).<br/>ok<br/>15&gt; X = 2+3.<br/>6<br/>16&gt;</pre>                                                                                                                                                                                                                                                                               |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Table of Contents                                                                                                                                                                                                                                                                                                                                                                                                                                            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Processes and Messages                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |                                                                                                                                           |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                             | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       | CS 418 – Jan. 6, 2016                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | 35 / 26                                                                                                                                   |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Table of Contents                                                                                                                                                                                                                                                                                                                                                                                                                                            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | Introduction to Erlang                                                                                                                    | CS 418 – Jan. 6, 2016                                                                                                                                                   | 36 / 26                                                                                                                    |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Table of Contents                                                                                                                                                                                                                                                                                                                                                                                                                                            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Processes and Messages                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |                                                                                                                                           | Objectives                                                                                                                                                              |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                             | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       | CS 418 – Jan. 6, 2016                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | 39 / 26                                                                                                                                   |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Table of Contents                                                                                                                                                                                                                                                                                                                                                                                                                                            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | Introduction to Erlang                                                                                                                    | CS 418 – Jan. 6, 2016                                                                                                                                                   | 40 / 26                                                                                                                    |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                             | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       | CS 418 – Jan. 6, 2016                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | 40 / 26                                                                                                                                   |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Table of Contents                                                                                                                                                                                                                                                                                                                                                                                                                                            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | Processes and Messages                                                                                                                    | CS 418 – Jan. 9, 2017                                                                                                                                                   | 1 / 20                                                                                                                     |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                             | Introduction to Erlang                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       | CS 418 – Jan. 6, 2016                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | 40 / 26                                                                                                                                   |                                                                                                                                                                         |                                                                                                                            |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |
| Table of Contents                                                                                                                                                                                                                                                                                                                                                                                                                                            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | Mark Greenstreet                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | Processes and Messages                                                                                                                    | CS 418 – Jan. 9, 2017                                                                                                                                                   | 2 / 20                                                                                                                     |                                                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |                       |         |

## Processes – Overview

- The built-in function `spawn` creates a new process.
- Each process has a process-id, pid.
  - The built-in function `self()` returns the pid of the calling process.
  - `spawn` returns the pid of the process that it creates.
  - The simplest form is `spawn(Fun)`.
    - A new process is created as child.
    - The pid of the child process is returned to the caller of `spawn`.
    - The function `Fun` is invoked with no arguments in that process.
    - The parent process and the child process are both running.
    - When `Fun` returns, the child process terminates.

## Processes – a friendly example

```
hello(N) ->
 [spawn(fun() -> io:format(
 "Hello world from process ~b~n", [I])
 end)
 || I <= lists:seq(1, N)
].
```

### Running the code:

```
>>> c(processes).
{ok,processes}
2>> procs:hello(3).
Hello world from process 1
Hello world from process 2
Hello world from process 3
[<0.40.0>,<0.41.0>,<0.42.0>]
```

## Reactive Processes and Tail Recursion

- Often we want processes that do more than add two numbers together.
- We want processes that wait, receive a message, process the message, and then wait for the next message.
- In Erlang, we do this with recursive functions for the child process:

```
accProc(Tally) ->
receive
 N when is_integer(N) -> 8>> BPid = proc:accumulator() .
 <0.53.0>
 accProc(Tally+N);
 {Pid, Total} -> 9>> BPid ! 2.
 Pid ! Tally,
 accProc(Tally);
 end.
end.
accumulator() ->
spawn(fun() -> 10>> BPid ! 3.
 accumulator();
 spawn(fun() -> 12>> receive Tl -> Tl end.
 accProc(0)
 end).
6
```

CS 418 – Jan. 9, 2017 7 / 20

## Reactive Processes and Tail Recursion

- Often, we want processes that do more than add two numbers together.
- We want processes that wait, receive a message, process the message, and then wait for the next message.
- In Erlang, we do this with recursive functions for the child process:

```
accProc(Tally) ->
receive
 N when is_integer(N) -> 13>> BPid ! 4.
 <0.53.0>
 accProc(Tally+N);
 {Pid, Total} -> 14>> BPid ! {self(), total},
 <0.33.0>, total>;
 Pid ! Tally,
 accProc(Tally);
 end.
15>> BPid ! 5.
5
16>> BPid ! 6.
6
17>> BPid ! {self(), total},
<0.33.0>, total>;
18>> receive T2 -> T2 end.
10
19>> receive T3 -> T3 end.
21
```

CS 418 – Jan. 9, 2017 7 / 20

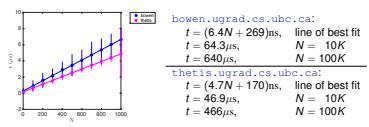
CS 418 – Jan. 9, 2017 7 / 20

## Timing Measurements

- We write parallel code to solve problems that would take too long on a single CPU.
- To understand performance trade-offs, I'll measure the time for some common operations in Erlang programs:
  - The time to make  $N$  recursive tail calls.
  - The time to spawn an Erlang process.
  - The time to send and receive messages:
    - Short messages.
    - Messages consisting of lists of varying lengths.

## Timing Measurements

### Tail Call Time



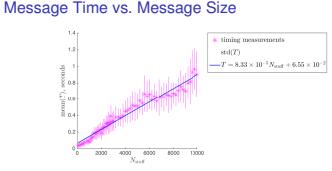
bowen.ugrad.cs.ubc.ca:  
 $t = 6.4N + 269\text{ns}$ , line of best fit  
 $t = 64.3\mu\text{s}$ ,  $N = 10K$

thetis.ugrad.cs.ubc.ca:  
 $t = 4.7N + 170\text{ns}$ , line of best fit  
 $t = 46.9\mu\text{s}$ ,  $N = 10K$

bowen.ugrad.cs.ubc.ca:  
 $t = 4.06N + 466\text{ns}$ ,  $N = 100K$

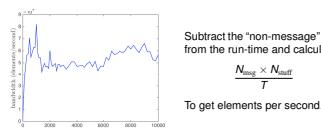
- Measurement: start the timing measurement, make  $N$  tail calls, and end the timing measurement.
- The measurements on this slide and throughout the lecture were made using the `timemeasure` function from the course Erlang library.
  - `timemeasure(Fun)` repeatedly calls `Fun` until about one second has elapsed. It then reports the average time and standard deviation.
  - `timemeasure` has lots of options.

## Message Time vs. Message Size



- Set-up: as on the previous slide. This time each message consists of a list of  $N_{\text{msg}}$  small integers.
- Each process sends and receives 5000 messages per run.
- The slope of the line divided by 5000 is the time per element:  
 $\sim 17\text{ns}/\text{message}$  on `thetis.ugrad.cs.ubc.ca, erts 18.2.`

## Bandwidth vs. Message Size



Subtract the “non-message” time from the run-time and calculate:

$$\frac{N_{\text{msg}} \times N_{\text{msg}}}{T}$$

To get elements per second.

- Bandwidth grows rapidly with message length for  $N_{\text{msg}} < 1000$ , then drops.
  - Short messages have low bandwidth due to fixed overheads with each message.
  - I'm guessing that bandwidth drops some for messages with more than 1000 elements because the Erlang runtime is somehow optimized for short messages.

- Processes are easy to create in Erlang.
  - The `spawn` mechanism can be used to start other processes on the same CPU or on machines spread around the internet.
- Processes communicate through messages
  - Message passing is asynchronous.
  - The receiver can use patterns to select a desired message.
- Reactive processes are implemented with tail-recursive functions.
- Interprocess operations are much slower than local ones
  - This is a key consideration in designing parallel programs.
  - We'll learn why when we look at parallel architectures later this month.

## Summary

### Preview

|                                             |                                                                                                                                                       |
|---------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------|
| January 11: Reduce                          | Reading: <a href="#">Learn You Some Erlang, Errors and Exceptions through A Short Visit to Common Data Structures</a>                                 |
| January 13: Scan                            | Reading: Lin & Snyder, chapter 5, pp. 112–125                                                                                                         |
|                                             | Mini-Assignment: <a href="#">Mini-Assignment 2 due 10:00am</a>                                                                                        |
| January 16: Generalized Reduce and Scan     | Homework: <a href="#">Homework 1 due 11:59pm</a> for early-bird bonus (11:59pm)<br><a href="#">Homework 2 goes out (due Feb. 1) – Reduce and Scan</a> |
| January 18: Reduce and Scan Examples        | Homework: <a href="#">Homework 1 due 11:59pm</a>                                                                                                      |
| January 20–27: Parallel Architecture        |                                                                                                                                                       |
| January 29–February 6: Parallel Performance |                                                                                                                                                       |
| February 8–17: Parallel Sorting             |                                                                                                                                                       |

- If your process is waiting for a message that never arrives, e.g. because
  - You misspelled a tag for a message, or
  - The receive pattern is slightly different than the message that was sent, or
  - Something went wrong in the sending process, and it died before sending the message, or
  - You got the message ordering slightly wrong, and there's a cycle of processes waiting for each other to send something, or
  - ...
- Then your process can wait forever, your Erlang shell can hang, and it's a very unhappy time in life.
  - See [Time Out in LVE](#).
  - Note: time-outs are great for debugging. They should be used with great caution elsewhere because they are sensitive to changes in hardware, changes in the scale of the system, and so on.
- Try it (e.g. with the example from slide 7).
- Don't forget to delete (or comment out) such debugging output before releasing your code.

CS 418 – Jan. 9, 2017 14 / 20

## Time Outs

- If your process is waiting for a message that never arrives, e.g. because
  - You misspelled a tag for a message, or
  - The receive pattern is slightly different than the message that was sent, or
  - Something went wrong in the sending process, and it died before sending the message, or
  - You got the message ordering slightly wrong, and there's a cycle of processes waiting for each other to send something, or
  - ...
- Then your process can wait forever, your Erlang shell can hang, and it's a very unhappy time in life.
  - See [Time Out in LVE](#).
  - Note: time-outs are great for debugging. They should be used with great caution elsewhere because they are sensitive to changes in hardware, changes in the scale of the system, and so on.

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9, 2017 22 / 20

CS 418 – Jan. 9,



## Generalize Reduce and Scan

Mark Greenstreet

CpSc 418 – Jan. 16, 2017

### Outline:

- Reduce in Erlang
- Scan in Erlang

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license (<http://creativecommons.org/licenses/by/4.0/>)

Mark Greenstreet Generalize Reduce and Scan CS 418 – Jan. 16, 2017 1 / 13

## Objectives

- Understand relationship between reduce and scan
  - ▶ Both are tree walks.
  - ▶ The initial combination of values from leaves is identical.
  - ▶ Reduce propagates the grand total down the tree.
  - ▶ Scan propagates the total ‘everything to the left’ down the tree.
- Generalized Reduce and Scan
  - ▶ Understand the role of the *Leaf*, *Combine*, and *Root* functions.
  - ▶ Understand the use of higher-order functions to implement reduce and scan.
  - The CS418 class library
    - ▶ Able to create a tree of processes.
    - ▶ Able to distribute data and tasks to those processes.
    - ▶ Able to use the *reduce* and *scan* functions from the library.
    - Know where to find more information.

## Reduce in Erlang

- Build a tree.
- Each process creates a lists of random digits.
- The processes meet at a barrier so we can measure the time to count the 3s.
- Each process counts its threes.
- The processes use reduce to compute the grand total.
- Each process reports the grand total and its own tally.
- The root process reports the time for the local tallies and the reduce.
- Get the code at

<http://www.ugrad.cs.ubc.ca/~cs418/2016-2/lecture/01-16/code/reduce.erl>

## The Reduce Pattern

- It's a parallel version of *fold*, e.g. `lists:foldl`.
- Reduce is described by three functions:
  - ▶ *Leaf()*: What to do at the leaves, e.g. `fun() -> count3s(Data)`.
  - ▶ *Combine()*: What to do at the root, e.g. `fun(Left, Right) -> Left+Right`.
  - ▶ *Root()*: What to do with the final result. For count 3s, this is just the identity function.

## The wtree module

- Part of the [course Erlang library](#).
- Operations on worker trees”
  - ▶ `wtree:create(NProc)` -> `[pid()]`. Create a list of *NProc* processes, organized as a tree.
  - ▶ `wtree:broadcast(W, Task, Arg)` -> `ok`. Execute the function *Task* on each process in *W*. Note: *W* means “worker pool”.
  - ▶ `wtree:reduce(P, Leaf, Combine, Root)` -> `term()`. A generalized reduce.
  - ▶ `wtree:reduce(P, Leaf, Combine)` -> `term()`. A generalized reduce where *Root* defaults to the identity function.

Mark Greenstreet Generalize Reduce and Scan CS 418 – Jan. 16, 2017 5 / 13

## Store Locally

- Communication is expensive – each process should store its own data whenever possible.
- How do we store data in a functional language?
  - ▶ Our processes are implemented as Erlang functions that receive messages, process the message, and make a tail-call to be ready to receive the next message.
  - ▶ We add a parameter to these functions, *State*, that is a mapping from *Keys* to *Values*.
- What this means when we write code:
  - ▶ Functions such as *Leaf* for *wtree:reduce* or *Task* for *wtree:broadcast* have a parameter for *State*.
  - ▶ `worker:put(State, Key, Value)` -> `NewState`. Create a new version of *State* that associates *Value* with *Key*.
  - ▶ `worker:put(State, Key, Default)` -> `Default`. Return the *Value* associated with *Key* in *State*. If no such value is found, *Default* is returned. Note: *Default* can be a function which case it is called to determine a default value – see the documentation.

## Count3s using wtree

```
count3s_per(N, P) ->
 W = wtree:create(P),
 wtree:list(W, N, 10, 'Data'),
 wtree:reduce(W,
 fun(ProcState) ->
 count3s_workers:get(ProcState, 'Data'))
 end,
 fun(Left, Right) -> Left+Right
end.
```

## Scan in Erlang

- Remarkably like reduce.
- Reduce has
  - ▶ an upward pass to compute the grand total
  - ▶ a downward pass to broadcast the grand total.
- Scan has
  - ▶ an upward pass where the grand total – **just like reduce**
  - ▶ On the downward pass, we compute the total of all elements to the left of each subtree.
- Get the code at

<http://www.ugrad.cs.ubc.ca/~cs418/2016-2/lecture/01-16/code/scan.erl>

Mark Greenstreet Generalize Reduce and Scan CS 418 – Jan. 16, 2017 6 / 13

## The Scan Pattern

- It's a parallel version of *mapfold*, e.g. `lists:mapfoldl` and `lists:mapfoldr`.
- `wtree:scan(Leaf1, Leaf2, Combine, Acc0)`
  - ▶ *Leaf1 (ProcState)* -> *Value*. Each worker process computes its *Value* based on its *ProcState*.
  - ▶ *Combine (Left, Right)* -> *Value*. Combine values from sub-trees.
  - ▶ *Leaf2 (ProcState, AccIn)* -> *ProcState*. Each worker updates its state using the *AccIn* value – i.e. the accumulated value of everything to the worker’s “left”.
  - ▶ *Acc0*: The value to use for *AccIn* for the leftmost nodes in the tree.

## Scan example: prefix sum

```
prefix_sum_per(W, Key1, Key2) ->
 wtree:scan(W,
 fun(ProcState) -> % Leaf
 lists:sum(wtree:get(ProcState, Key1)),
 fun(ProcState, AccIn) -> % Leaf2
 wtree:put(ProcState, Key2,
 prefix_sum(wtree:get(ProcState, Key1), AccIn))
 end,
 fun(Left, Right) -> % Combine
 Left + Right
 end,
 0 % Acc0
).
```

```
prefix_sum(L, Acc0) ->
 element(1,
 lists:mapfoldl(fun(X, Y) -> Sum = X+Y, {Sum, Sum}, end,
 Acc0, L)).
```

Mark Greenstreet Generalize Reduce and Scan CS 418 – Jan. 16, 2017 7 / 13

## Scan

Mark Greenstreet  
CpSc 418 – Jan. 20, 2016

## Preview

### January 18: Reduce and Scan Examples

Homework: Homework 1 due 11:59pm

January 20: Finish Reduce and Scan

Mini-assessments: Mini assignment 3 goes out.

January 23: Architecture Review

Reading: Parallel Chapter 2, through section 2.2

January 27: Shared Memory Architectures

Reading: Pacheco, Chapter 2, through section 2.3

Mini-assessments: Mini assignment 3 due, 10am.

January 27: Message Passing

January 30: HW 2 (earlybird due (11:59pm), HW 3 goes out.

February 1: HW 2 due (11:59pm).

February 8-17: Parallel Sorting

February 8: HW 1 due (11:59pm).

February 17: HW 3 due (11:59pm).

February 27: TBD

March 1: Midterm

Mark Greenstreet Generalize Reduce and Scan CS 418 – Jan. 16, 2017 10 / 13

## Parallel Prefix Sum

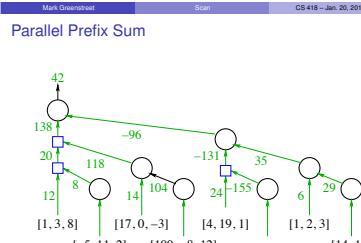
Mark Greenstreet

Mark Greenstreet Generalize Reduce and Scan CS 418 – Jan. 20, 2016 1 / 15

## Parallel Prefix Sum

Mark Greenstreet

Mark Greenstreet Generalize Reduce and Scan CS 418 – Jan. 20, 2016 4 / 15



Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

## Scan example: prefix sum

```
prefix_sum_per(W, Key1, Key2) ->
 wtree:scan(W,
 fun(ProcState) -> % Leaf
 lists:sum(wtree:get(ProcState, Key1)),
 fun(ProcState, AccIn) -> % Leaf2
 wtree:put(ProcState, Key2,
 prefix_sum(wtree:get(ProcState, Key1), AccIn))
 end,
 fun(Left, Right) -> % Combine
 Left + Right
 end,
 0 % Acc0
).
```

```
prefix_sum(L, Acc0) ->
 element(1,
 lists:mapfoldl(fun(X, Y) -> Sum = X+Y, {Sum, Sum}, end,
 Acc0, L)).
```

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

## Scan

Mark Greenstreet

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

## Parallel Prefix Sum

Mark Greenstreet

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

## Parallel Prefix Sum

Mark Greenstreet

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

## Parallel Prefix Sum

Mark Greenstreet

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

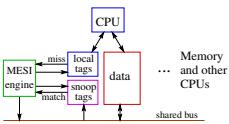
Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 2016 4 / 15

Mark Greenstreet Scan CS 418 – Jan. 20, 201



## Implementing MESI: Snooping



- Caches read and write main memory over a shared memory bus.
- Each cache has two copies of the tags: one for the CPU, the other for the bus.
- If the cache sees another CPU reading or writing a block that is in its cache, it takes the action specified by the MESI protocol.

## Implementing MESI: Directories

- Main memory keeps a copy of the data and
  - a bit-vector that records which processors have copies, and
  - a bit to indicate that one processor has a copy and it may be modified.
- A processor accesses main memory as required by the MESI protocol.
  - The memory unit sends messages to the other CPUs to direct them to take actions as defined by the protocol.
  - The ordering of these messages ensures that memory stays consistent.
- Comparison:
  - Snooping is simple for machines with a small number of processors.
  - Directory methods scale better to large numbers of processors.

## Sequential Consistency

Memory is said to be **sequentially consistent** if

- All memory reads and writes from all processors can be arranged into a single, sequential order, such that:
  - The operations for each processor occur in the global ordering in the same order as they did on the processor.
  - Every write gets the value of the preceding write to the same address.
- Sequential consistency corresponds to what programmers think "ought" to happen.
  - Very similar to "serializability" for database transactions.
- MESI guarantees sequential consistency

## Coding Break

### Weak Consistency

- CPUs typically have "write-buffers" because memory writes often come in bursts.
- Typically, reads can move ahead of writes to maximize program performance.
- Why?
  - Because there may be instructions waiting for the data from a load.
  - A transition from "shared" to "modified" requires **all** processors – this can take a long time.
  - Memory writes don't happen until the instruction commits.
- This means that real computers don't guarantee sequential consistency.
  - Warning:** classical algorithms for locks and shared buffers fail when run on a real machine!

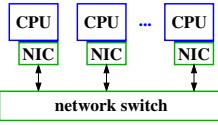
### Preview

| January 27: Distributed-Memory Machines     |                                                              |
|---------------------------------------------|--------------------------------------------------------------|
| Reading:                                    | Pacheco, Chapter 2, Sections 2.4 and 2.5.                    |
| Mini Assignments:                           | Mini 4 goes out.                                             |
| January 30: Parallel Performance: Speed-up  | Reading: Pacheco, Chapter 2, Section 2.6.                    |
| Homework:                                   | HW 2 earlybird (11:59pm). HW 3 goes out.                     |
| February 1: Parallel Performance: Overheads | Homework: HW 2 earlybird (11:59pm).                          |
| February 4: Parallel Performance: Wrap-up   | Homework: HW 3 due (10am).                                   |
| February 6: Parallel Performance: Wrap Up   | Homework: HW 3 due (11:59pm).                                |
| January 8–February 15: Parallel Sorting     | Homework (Feb. 15): HW 3 earlybird (11:59pm), HW 4 goes out. |
| February 17: Map-Reduce                     | Homework: HW 3 due (11:59pm).                                |
| February 27: TBD                            |                                                              |
| March 1: Midterm                            |                                                              |

### Cache Design Trade-Offs (2 of 2)

- Associativity:**
  - Increasing associativity generally reduces the number of conflict misses.
  - Increasing associativity makes the cache hardware more complicated.
  - Typical caches are direct mapped to four- or eight-way associative.
  - Associativity doesn't need to be a power of two!
- Other stuff**
  - cache inclusion: is everything in the L1 also in the L2?
  - interaction with virtual memory: are cache addresses virtual or physical?
  - coherence protocol details:
    - Example: Intel uses MESIF, the "F" stands for "forwarding". If a processor has a read miss, and another cache has a copy, one of the caches with a copy will be the "forwarding cache". The forwarding cache provides the data because it's much faster than main memory.
    - error detection and creation – caches + cosmic rays = flipped bits.
    - and all kinds of other optimizations that are beyond the scope of this class.

### Message Passing Computers

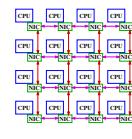


- Multiple CPUs
- Communication through a network:
  - Commodity networks for small clusters.
  - Special high-performance networks for super-computers
- Programming model:
  - Explicit message passing between processes (like Erlang)
  - No shared memory or variables.

### Network Topologies

- Network topologies are to the message-passing community what cache-coherence protocols are to the shared-memory people:
  - Lots of papers have been published.
  - Machine designers are always looking for better networks.
  - Network topology has a strong impact on performance, the programming model, and the cost of building the machine.
- A message-passing machine may have multiple networks:
  - A general purpose network for sending messages between machines.
  - Dedicated networks for reduce, scan, and synchronization:
    - The reduce and scan networks can include ALUs (integer and/or floating point) to perform common operations such as sums, max, product, all, any, etc. in the networking hardware.
    - A synchronization network only needs to carry a few bits and can be designed to minimize latency.

### Meshes



- Advantages:**
  - Easy to implement; chips and circuit boards are effectively two-dimensional.
  - Cross section bandwidth grow with number of processors – more specifically, bandwidth grows as  $\sqrt{P}$ .
- Disadvantages:**
  - Worst-case latency grows as  $\sqrt{P}$ .
  - Edges of mesh are "special cases."

### Implementing MESI: Directories

- Main memory keeps a copy of the data and
  - a bit-vector that records which processors have copies, and
  - a bit to indicate that one processor has a copy and it may be modified.
- A processor accesses main memory as required by the MESI protocol.
  - The memory unit sends messages to the other CPUs to direct them to take actions as defined by the protocol.
  - The ordering of these messages ensures that memory stays consistent.
- Comparison:
  - Snooping is simple for machines with a small number of processors.
  - Directory methods scale better to large numbers of processors.

## Sequential Consistency

Memory is said to be **sequentially consistent** if

- All memory reads and writes from all processors can be arranged into a single, sequential order, such that:
  - The operations for each processor occur in the global ordering in the same order as they did on the processor.
  - Every write gets the value of the preceding write to the same address.
- Sequential consistency corresponds to what programmers think "ought" to happen.
  - Very similar to "serializability" for database transactions.
- MESI guarantees sequential consistency

## Coding Break

### Shared Memory and Performance

- Shared memory can offer better performance than message passing because:
  - High bandwidth: the buses that connect the caches can be very wide, especially if the caches are on a single chip.
  - Low latency: the hardware handles moving the data – no operating system calls and context-switch overheads.
- But, shared memory doesn't scale as well as message passing
  - For large machines, the latency of directory accesses can severely degrade performance.
    - In a message passing machine, each CPU has its own memory, nearby and fast.
    - For shared memory, each CPU has part of the shared main memory – accessing a directory may require accessing the memory of a distant CPU.
  - Shared memory moves the data after the cache miss
    - This stalls threads.
    - message passing can send data in advance and avoid these stalls

### Summary

- Shared Memory Architectures**
  - Use cache-coherence protocols to allow each processor to have its own cache while maintaining (almost) the appearance of having one shared memory for all processors.
    - A typical protocol: MESI
    - The protocol can be implemented by snooping or directories.
  - Using cache-memory interconnect for interprocessor communication provides:
    - High-bandwidth
    - Low-latency, but watch out for fences, etc.
    - High cost for large scale machines.
- Shared-Memory Programming**
  - Need to avoid interference between threads.
    - Assertion reasoning (e.g. invariants) are crucial, much more so than in sequential programming.
    - There are too many possible interleavings to handle intuitively.
    - In practice, we don't formally prove complete programs, but we use the ideas of formal reasoning.
  - Real computers don't provide sequential consistency.
    - Use a thread library.

### Cache Design Trade-Offs (1 of 2)

- Capacity:** Larger caches have lower miss rates, but longer access times. This motivates using multiple levels of caches.
  - L1: closest to the CPU, smallest capacity (16-64Kbytes), fastest access (1-3 clock cycles).
  - L2: typically 128bytes to 1MByte, 5-10 cycle access time.
  - L3: becoming common, several Mbytes of capacity.
- Block Size:**
  - Larger blocks can lower miss rate by exploiting spatial locality.
  - Larger blocks can raise miss rate due to conflict and coherence.
  - Larger blocks increase miss penalty by requiring more time to transfer all that data.
  - Typical block sizes are 16 to 256 bytes – sometimes block size changes with cache level.

### Classifying Cache Misses

- Compulsory:** The first reference to a cache block will cause a miss.
  - Note that the first access should be a write – otherwise the location is uninitialized.
  - A cache can avoid stalling the processor by using "allocate on write".
  - If a miss is a write, assign a block for the line, start the main memory read, track which bytes have been written, and merge with the data from memory when it arrives.
- Capacity:** The cache is not big enough to hold all of the data used by the program.
- Conflict:** Many active memory locations map to the same cache index.
  - If there are more such references than the associativity of the cache, these will cause conflict misses.
- Coherence:** A cache block was evicted because another CPU was writing to it.
  - A subsequent read incurs a cache miss.

### Objectives

#### Message Passing Computers

##### Mark Greenstreet

CpSc 418 – Jan. 27, 2017

##### Outline:

- Network Topologies
- Performance Considerations
- Examples

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license  
<http://creativecommons.org/licenses/by/4.0/>

### False Sharing

- False sharing occurs when two CPUs are actively writing different words in the same cache block.
  - Each write forces the other CPU to invalidate its cache block.
  - Each read forces the other CPU to change its cache block from **modified** or **exclusive** to **shared**.
- Example: count 3s
  - Here's an implementation with awful performance.
  - We create a global array of ints to hold the accumulators for each process.
  - Each time a process finds a 3, it writes to its element in the array.
  - This forces the other CPUs whose accumulators are in the same block to invalidate their cache entry.
  - This turns accumulator accesses into main memory accesses.
  - And these accesses are serialized: one CPU at a time.

### Message Passing Computers

##### Mark Greenstreet

CpSc 418 – Jan. 27, 2017

##### Outline:

- Network Topologies
- Performance Considerations
- Examples

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license  
<http://creativecommons.org/licenses/by/4.0/>

### Some simple message-passing clusters

- 25 linux workstations (e.g. lin01 ... lin25.ugrad.cs.ubc.ca) and standard network routers.
  - A good platform for learning to use a message-passing cluster.
  - But, we'll figure out that network bandwidth and latency are key bottlenecks.
- A "blade" based cluster, for example:
  - 16 "blades" each with 4 6-core CPU chips, and 32G of DRAM.
  - An "infinband" or similar router for about 10-100 times the bandwidth of typical ethernet.
  - The price tag is ~\$300K.
    - Great if you need the compute power.
    - But, we won't be using one in this class.

### The Sunway TaihuLight

##### Mark Greenstreet

Cs 418 – Jan. 27, 2017

- The world's fastest (Linpack) super-computer (as of June 2016)
  - 40,960 multicore CPUs
    - 256 cores per CPU chip.
    - 1.45GHz clock frequency, 8 flops/core/cycle.
  - Total of 10,485,760 cores
  - LINPACK performance: 93 Pflops
  - Power consumption 15MW (computer) + cooling (unspecified)
  - Tree-like
    - Five levels of hierarchy.
    - Each level has a high-bandwidth switch.
    - Some levels (all?) are fully-connected for that level.
  - Programming model: A version linux with MPI tuned for this machine.
  - For more information, see [Report on the Sunway TaihuLight System](#), J. Dongarra, June 2016.

### The Westgrid Clusters

##### Mark Greenstreet

Cs 418 – Jan. 27, 2017

- Clusters at various Western Canadian Universities (including UBC).
  - Up to 9600 cores.
  - Available for research use.

### Ring-Networks

- Advantages: simple.
- Disadvantages:
  - Worst-case latency grows as  $O(P)$  where  $P$  is the number of processors.
  - Easily congested – limited bandwidth.

### Star Networks

##### Mark Greenstreet

Cs 418 – Jan. 27, 2017

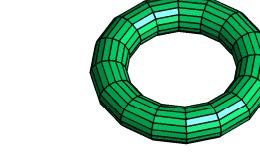
- Advantages:
  - Low-latency – single hop between any two nodes
  - High-bandwidth – no contention for connections with different sources and destinations.
- Disadvantages:
  - Amount of routing hardware grows as  $O(P^2)$ .
  - Requires lots of wires, to and from switch – Imagine trying to build a switch that connects to 1000 nodes!
- Summary
  - Surprisingly practical for 10-50 ports.
  - Hierarchies of cross-bars are often used for larger networks.

### A crossbar switch

##### Mark Greenstreet

Cs 418 – Jan. 27, 2017

### Tori



- Advantages:
  - Has the good features of a mesh, and
  - No special cases at the edges.
- Disadvantages:
  - Worst-case latency grows as  $\sqrt{P}$ .

### Hypercubes

##### Mark Greenstreet

Cs 418 – Jan. 27, 2017

#### A 0-dimensional (1 node), radix-2 hypercube

##### Mark Greenstreet

Cs 418 – Jan. 27, 2017

### Hypercubes

##### Mark Greenstreet

Cs 418 – Jan. 27, 2017

#### A 1-dimensional (2 node), radix-2 hypercube

##### Mark Greenstreet

Cs 418 – Jan. 27, 2017

## Hypercubes

A 2-dimensional (4 node), radix-2 hypercube



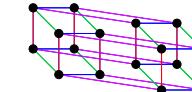
## Hypercubes

A 3-dimensional (8 node), radix-2 hypercube



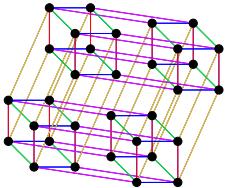
## Hypercubes

A 4-dimensional (16 node), radix-2 hypercube



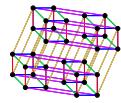
## Hypercubes

A 5-dimensional (32 node), radix-2 hypercube



## Hypercubes

A 5-dimensional (32 node), radix-2 hypercube



## Dimension Routing

```
% Send a message, msg, from node src to node dst
for i = 1:d % d is dimension of the hypercube
 if(bit(i, src) != bit(i, dst)) % if different for dimension i
 send(msg, link[i]); % then send msg to our i-neighbour
 end
end
```

- Advantages**
  - Small diameter ( $\log N$ )
  - Lots of bandwidth
  - Easy to partition.
  - Simple model for algorithm design.
- Disadvantages**
  - Needs to be squeezed into a three-dimensional universe.
  - Lots of long wires to connect nodes.
  - Design of a node depends on the size of the machine.

## Performance Considerations

- Bandwidth**
  - How many bytes-per-second can we send between two processors?
  - May depend on which two processors: neighbours may have faster links than spanning the whole machine
  - Bisection bandwidth: find the **worst** way to divide the processors into sets of  $P/2$  processors each.
    - How many bytes-per-second can we send between the two partitions?
    - If we divide this by the number of processors, we typically get a much smaller value that the peak between the two processors.
- Latency**
  - How long does it take to send a message from one processor to another?
  - Typically matters the most for short messages.
- Cost**
  - How expensive is the interconnect – it may dominate the total machine cost.
  - Cost of the network interface hardware.
  - Cost of the cables.

## Real-life networks

- InfiniBand is becoming increasingly prevalent
- Peak bandwidths > 6GB/sec.
- achieved bandwidths of 2-3GB/s.
- Support for RDMA and “one-sided” communication
- CPU A can read or write a block of memory residing with CPU B.
- Often, networks include trees for synchronization (e.g. barriers), and common reduce and scan operations.
- The MPI (message-passing interface) evolves to track the capabilities of the hardware.

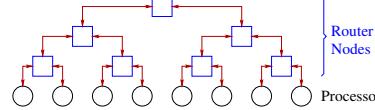
## Summary

- Message passing machines have an architecture that corresponds to the message-passing programming paradigm.
- Message passing machines can range from
  - Clusters of PCs with a commodity switch.
  - Clouds: lots of computers with a general purpose network.
  - Super-computers: lots of compute nodes tightly connected with high-performance interconnect.
- Many network topologies have been proposed:
  - Performance and cost are often dominated by network bandwidth and latency.
  - The network can be more expensive than the CPUs.
  - Peta-flops or other instruction counting measures are an indirect measure of performance.
- Implications for programmers
  - Location matters
  - Communication costs of algorithms is very important
  - Heterogeneous computing is likely in your future.

## Preview

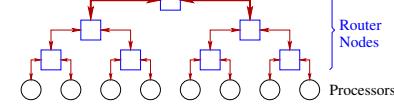
- January 30: Parallel Performance: Speed-up  
Reading: Pachico, Chapter 2, Section 2.6.  
Homework: HW 2 earlybird (11:59pm), HW 3 goes out.
- February 1: Parallel Performance: Overheads  
Reading: Pachico, Chapter 2, Section 2.7.
- February 3: Parallel Performance: Models  
Mini Assignments Mini 4 due (10am)
- February 6: Parallel Performance: Wrap Up  
January 6-February 15: HW 3 earlybird (11:59pm), HW 4 goes out.
- February 17: Map-Reduce  
Homework: HW 3 due (11:59pm).
- February 27: TBD
- March 1: Midterm
- There will be readings assigned from [Programming Massively Parallel Processors](#) starting after the midterm.
  - Make sure you have a copy. Note: this year, we'll make sure the course works with either the 2<sup>nd</sup> or 3<sup>rd</sup> edition.

## Trees



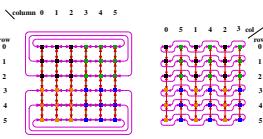
- Simple network: number of routing nodes = number of processors – 1.
- Wiring:  $O(\log N)$  extra height ( $O(N \log N)$ ) extra area.
- Wiring:  $O(\sqrt{N} \log N)$  extra area for H-tree.
- Low-latency:  $O(\log N)$  + wire delay.
- Low-bandwidth: bottleneck at root.

## Fat-Trees



- Use  $M^{\alpha}$  parallel links to connect subtrees with  $M$  leaves.
- $0 \leq \alpha \leq 1$ 
  - $\alpha = 0$ : simple tree
  - $\alpha = 1$ : strange crossbar
- Fat-trees are “universal”
  - For  $\frac{2}{3} < \alpha < 1$  a fat-tree interconnect with volume  $V$  can simulate **any** interconnect that occupies the same volume with a time overhead that is poly-log factor of  $N$ .

## From a mesh to a torus (1/2)



- Fold left-to-right, and make connections where the left and right edges meet.
- Now, we've got a cylinder.
- Note that there are no “long” horizontal wires: the longest wires jump across one processor.

## Review

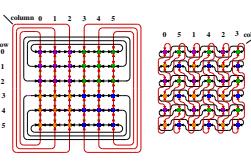
## Review

- Consider a machine with 4096 processors.
- What is the maximum latency for sending a message between two processors (measured in network hops)?
- Is it a ring?
- Is it a crossbar?
- Is it a 2-D mesh?
- Is it a 3-D mesh?
- Is it a hypercube?
- Is it a binary tree?
- Is it a radix-4 tree?

## Supplementary Material

- [Message-passing origami](#): how to fold a mesh into a torus.
- [How big is a hypercube](#): it's all about the wires.

## From a mesh to a torus (2/2)



- Fold top-to-bottom, and make connections where the top and bottom edges meet.
- Now, we've got a torus.
- Again there are no “long” wires.

## How big is a hypercube?

- Consider a hypercube with  $N = 2^d$  nodes.
- Assume each link can transfer one message in each direction in one time unit. The analysis here easily generalizes for links of higher or lower bandwidths.
- Let each node send a message to each of the other nodes.
- Using dimension routing.
  - Each node will send  $N/2$  messages for each of the  $d$  dimensions.
  - This takes time  $N/2$ .
  - As soon as one batch of messages finishes the dimension-0 route, that batch can continue with the dimension-1 route, and the next batch can start the dimension 0 route.
  - So, we can route with a throughput of  $\binom{N}{2}$  messages per  $N/2$  time.

## How big is a hypercube?

- Consider a hypercube with  $N = 2^d$  nodes.
- Assume each link can transfer one message in each direction in one time unit. The analysis here easily generalizes for links of higher or lower bandwidths.
- Let each node send a message to each of the other nodes.
- Using dimension routing,
  - we can route with a throughput of  $\binom{N}{2}$  messages per  $N/2$  time.
- Consider any plane such that  $N/2$  nodes are on each side of the plane.
  - $\frac{1}{2} \binom{N}{2}$  messages must cross this plane in  $N/2$  time.
  - This means that at least  $N - 1$  links must cross the plane.
  - The plane has area  $O(N)$ .

## How big is a hypercube?

- Consider a hypercube with  $N = 2^d$  nodes.
- Assume each link can transfer one message in each direction in one time unit. The analysis here easily generalizes for links of higher or lower bandwidths.
- Let each node send a message to each of the other nodes.
- Using dimension routing,
- we can route with a throughput of  $\binom{N}{2}$  messages per  $N/2$  time.
- Consider any plane such that  $N/2$  nodes are on each side of the plane.
  - The plane has area  $O(N)$ .
- Because the argument applies for any plane, we conclude that the hypercube has diameter  $O(\sqrt{N})$  and thus volume  $O(N^{\frac{d}{2}})$ .
- Asymptotically, the hypercube is all wire.

## Speed-Up

### Mark Greenstreet

CpSc 418 – Jan. 30, 2017

## Outline:

- [Measuring Performance](#)
- [Speed-Up](#)
- [Amdahl's Law](#)
- [The law of modest returns](#)
- [Superlinear speed-up](#)
- [Embarrassingly parallel problems](#)

## But first, USRA

### Summer Undergraduate Research Opportunities

- [Natural Sciences and Engineering Research Council \(NSERC\) Undergraduate Student Research Awards \(USRA\)](#)
  - Same process to apply for Science Undergraduate Research Experience (SURE) and Work Learn International Undergraduate Research Awards
- See what academic research really looks like
- Many research areas: ...
  - Google “usra cs usra” for full list of projects seeking students
- I have several project proposals:
  - Code for a fleet of smart wheelchairs for older adults
  - Numerical software for demonstrating correctness of robots and cyber-physical systems
  - 16 weeks, flexible schedule
  - You get paid!
  - Email potential sponsor ASAP (full applications due by Feb 10)

## Measuring Performance

- The main motivation for parallel programming is performance
  - Time: make a program run faster.
  - Space: allow a program to run with more memory.
- To make a program run faster, we need to know how fast it is running.
- There are many possible measures:
  - Latency: time from starting a task until it completes.
  - Throughput: the rate at which tasks are completed.
  - Key observation:

$$\text{throughput} = \frac{1}{\text{latency}} \cdot \text{sequential programming}$$

$$\text{throughput} \geq \frac{1}{\text{latency}} \cdot \text{parallel programming}$$

## Speed-Up

Simple definition:

$$\text{speed-up} = \frac{\text{time(sequential\_execution)}}{\text{time(parallel\_execution)}}$$

We can also describe speed-up as how many percent faster:

$$\% \text{faster} = (\text{speed-up} - 1) * 100\%$$

But beware of the spin:

- Is “time” latency or throughput?
- How big is the problem?
- What is the sequential version:
  - The parallel code run on one processor?
  - The fastest possible sequential implementation?
  - Something else?

More practically, how do we measure time?

## Speed-Up – Example

- Let's say that counting 3s of a million items takes 10ms on a single processor.
- If I run count 3s with four processes on a four CPU machine, and it takes 3.2ms, what is the speed-up?
- If I run count 3s with 16 processes on a four CPU machine, and it takes 1.8ms, what is the speed-up?
- If I run count 3s with 128 processes on a 32 CPU machine, and it takes 0.28ms, what is the speed-up?

## Time complexity

- What is the time complexity of sorting?
  - What are you counting?
  - Why do you care?
- What is the time complexity of matrix multiplication?
  - What are you counting?
  - Why do you care?

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license  
<http://creativecommons.org/licenses/by/4.0/>

January 2017

Ian M. Mitchell – USC Computer Science

1

CS 418 – Jan. 30, 2017

Speed-Up

2 / 20

Mark Greenstreet

Speed-Up

3 / 20

Mark Greenstreet

Speed-Up

4 / 20

Mark Greenstreet

Speed-Up

CS 418 – Jan. 30, 2017

5 / 20

Mark Greenstreet

Speed-Up

6 / 20

Mark Greenstreet

Speed-Up

7 / 20

## Big-O and Wall-Clock Time

- In our algorithms classes, we count "operations" because we have some belief that they have something to do with how long the actual program will take to execute.
- Or maybe not. Some would argue that we count "operations" because it allows us to use nifty techniques from discrete math.
- I'll take the position that the discrete math is nifty because it tells us something useful about what our software will do.
- In our architecture classes, we got the formula:
$$\text{time} = \frac{(\# \text{inst. executed}) * (\text{cycles/instruction})}{\text{clock frequency}}$$

- The approach in algorithms class of counting comparisons or multiplications, etc., is based on the idea that everything else is done in proportion to these operations.
- BUT**, in parallel programming, we can find that a communication between processes can take 1000 times longer than a comparison or multiplication.
- This may not matter if you're willing to ignore "constant factors."
- In practice, factors of 1000 are too big to ignore.

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 8 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 9 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 10 / 20

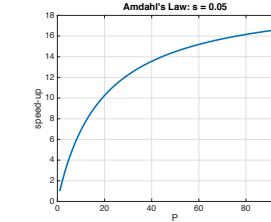
Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 11 / 20

## Amdahl's Law, 49 years later

## Amdahl's Law

- Given a sequential program where
  - fraction  $s$  of the execution time is inherently sequential.
  - fraction  $1-s$  of the execution time benefits perfectly from speed-up.
- The run-time on  $P$  processors is:
 
$$T_{\text{parallel}} = T_{\text{sequential}} * (s + \frac{1-s}{P})$$
- Consequences:
  - Define
 
$$\text{speed\_up} = \frac{T_{\text{sequential}}}{T_{\text{parallel}}}$$
  - Speed-up on  $P$  processors is at most  $\frac{1}{s}$ .
  - Gene Amdahl argued in 1967 that this limit means that parallel computers are only useful for a few special applications where  $s$  is very small.

## Amdahl's Law



## Amdahl's Law, 49 years later

Amdahl's law is not a physical law.

- Amdahl's law is mathematical theorem:
  - If  $T_{\text{parallel}}$  is  $(s + \frac{1-s}{P}) T_{\text{sequential}}$
  - and  $\text{speed\_up} = \frac{T_{\text{sequential}}}{T_{\text{parallel}}}$ ,
  - then for  $0 < s < 1$ ,  $\text{speed\_up} \leq \frac{1}{s}$ .
- Amdahl's law is also an **economic** law:
  - Amdahl's law was formulated when CPUs were expensive.
  - Today, CPUs are cheap
    - The cost of fabricating eight cores on a die is very little more than the cost of fabricating one.
    - Computer cost is dominated by the rest of the system: memory, disk, network, monitor, ...
- Amdahl's law assumes a fixed problem size.

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 8 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 9 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 10 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 11 / 20

## Amdahl's Law, 49 years later

## Amdahl's Law, 49 years later

- Amdahl's law is an **economic** law, not a **physical** law.
  - Amdahl's law was formulated when CPUs were expensive.
  - Today, CPUs are cheap (see previous slide)
- Amdahl's law assumes a fixed problem size
  - Many computations have  $s$  (sequential fraction) that decreases as  $N$  (problem size) increases.
  - Having lots of cheap CPUs available will
    - Change our ideas of what computations are easy and which are hard.
    - Determine what the "killer-apps" will be in the next ten years.
    - Ten years from now, people will just take it for granted that most new computer applications will be parallel.
  - Examples: see next slide

## Amdahl's Law, one more try



- We can have problems where the parallel work grows faster than the sequential part.
- Example: parallel work grows as  $N^{3/2}$  and the sequential part grows as  $\log P$ .

## The Law of Modest Returns

More bad news: ☺

- Let's say we have an algorithm with a sequential run-time  $T = (12m)N^4$ .
  - If we're willing to wait for one hour for it to run, what's the largest value of  $N$  we can use?
  - If we have 1000 machines, and perfect speed-up (i.e.,  $\text{speed\_up} = 1000$ ), now what is the largest value of  $N$  we can use?
  - What if the run-time is  $(5m)1.2^P$ ?
- The law of modest returns
  - Parallelism offers modest returns, unless the problem is of fairly low complexity.
  - Sometimes, modest returns are good enough: weather forecasting, climate models.
  - Sometimes, problems have huge  $N$  and low complexity: data mining, graphics, machine learning.

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 8 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 9 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 10 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 11 / 20

## Super-Linear Speed-up

Sometimes,  $\text{speed\_up} > P$ . ☺

- How does this happen?
  - Impossibility "proof": just simulate the  $P$  parallel processors with one processor, time-sharing  $P$  ways.
- Memory: a common explanation
  - $P$  machines have more main memory (DRAM)
  - and more cache memory and registers (total)
  - and more I/O bandwidth, ...
- Multi-threading: another common explanation
  - The sequential algorithm underutilizes the parallel capabilities of the CPU.
  - A parallel algorithm can make better use.
- Algorithmic advantages: once in a while, you win!
  - Simulation as described above has overhead.
  - If the problem is naturally parallel, the parallel version can be more efficient.
- BUT**: be very skeptical of super-linear claims, especially if  $\text{speed\_up} \gg P$ .

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 12 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 13 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 14 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 15 / 20

## Review Questions

- What is speed-up? Give an intuitive, English answer and a mathematical formula.
- Why can it be difficult to determine the sequential time for a program when measuring speed-up?
- What is Amdahl's law? Give a mathematical formula. Why is Amdahl's law a concern when developing parallel applications? Why in many cases is it not a show-stopper?
- Is parallelism an effective solution to problems with high big-O complexity? Why or why not?
- What is super-linear speed-up? Describe two causes.
- What is an embarrassingly parallel problem. Give an example.

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 16 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 17 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 18 / 20

Mark Greenstreet Speed-Up CS 418 – Jan. 30, 2017 19 / 20

## Embarrassingly Parallel Problems

## Embarrassingly Parallel Problems

- Problems that can be solved by a large number of processors with very little communication or coordination.
  - Rendering images for computer-animation: each frame is independent of all the others.
  - Brute-force searches for cryptography.
  - Analyzing large collections of images: astronomy surveys, facial recognition.
  - Monte-Carlo simulations: same model, run with different random values.
  - Don't be ashamed if your code is embarrassingly parallel!**
    - Embarrassingly parallel problems are great: you can get excellent performance without heroic efforts.
    - The only thing to be embarrassed about is if you don't take advantage of easy parallelism when it's available.

## Lecture Summary

### Parallel Performance

- Speed-up: [slide 5](#)
- Limits
  - Amdahl's Law, [slide 9](#)
  - Modest gains, [slide 15](#).
- Sometimes, we win
  - Super-linear speedup, [slide 16](#)
  - Embarrassingly Parallel Problems, [slide 17](#).

### Objectives

- Learn about main causes of performance loss:
  - Overhead
  - Non-parallelizable code
  - Idle processors
  - Resource contention
- See how these arise in message-passing, and shared-memory code.

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license <http://creativecommons.org/licenses/by/4.0/>.

### Performance-Loss

Mark Greenstreet

CpSc 418 – Feb. 1, 2017

### Outline:

- Overhead:** work the parallel code has to do that the sequential version avoids.
  - Communication and Synchronization
  - Extra computation, extra memory
- Limited parallelism**
  - Code that is inherently sequential or has limited parallelism
  - Idle processors
  - Resource contention

### Causes of Performance Loss

- Ideally, we would like a parallel program to run  $P$  times faster than the sequential version when run on  $P$  processors.
- In practice, this rarely happens because of:
  - Overhead**: work that the parallel program has to do that isn't needed in the sequential program.
  - Non-parallelizable code**: something that has to be done sequentially.
  - Idle processors**: There's work to do, but some processor are waiting for something before they can work on it.
  - Resource contention**: Too many processors overloading a limited resource.

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 1 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 2 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 3 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 4 / 22

## Communication Overhead

## Communication Overhead

- In a parallel program, data must be sent between processors.
- This isn't a part of the sequential program.
- The time to send and receive data is overhead.
- Communication overhead occurs with both shared-memory and message-passing machines and programs.
- Example: Reduce (e.g., Count 3s):
  - Communication between processes adds time to execution.
  - The sequential program doesn't have this overhead.

## Communication with shared-memory

- In a shared memory architecture:
  - Each core has its own cache.
  - The caches communicate to make sure that all references from different cores to the same address look like there is one, common memory.
  - It takes longer to access data from a remote cache than from the local cache. This creates overhead.
- False sharing** can create communication overhead even when there is no logical sharing of data.
  - This occurs if two processors repeatedly modify different locations on the same cache line.

## Communication overhead: example

- The [Principles of Parallel Programming](#) book considered an example of Count 3s (in C, with threads), where there was a global array, `int count[P]` where  $P$  is the number of threads.
  - Each thread (e.g. thread  $i$ ) initially sets its count, `count[i] = 0.`
  - Each time a thread encounters a 3, it increments its element in the array.
- The parallel version ran much slower than the sequential one.
  - Cache lines are much bigger than a single `int`. Thus, many entries for the `count` array are on the same cache line.
  - A processor has to get exclusive access to update the count for its thread.
  - This invalidates the copies held by the other processors.
  - This produces lots of cache misses and a slow execution.
- A better solution:
  - Each thread has a local variable for its count.
  - Each thread counts its threads using this local variable and copies its final total to the entry in the global array.

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 5 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 6 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 7 / 22

## Communication overhead with message passing

## Synchronization Overhead

- Parallel processes must coordinate their operations.
  - Access to shared data structures.
  - Writing to a file.
- For shared-memory programs (e.g., `pthreads` or `Java threads`), there are explicit locks or other synchronization mechanisms.
- For message passing (e.g., `Erlang` or `MPI`), synchronization is accomplished by communication.

## Computation Overhead

- A parallel program may perform computation that is not done by the sequential program.
  - Redundant computation: it's faster to recompute the same thing on each processor than to broadcast.
  - Algorithm: sometimes the fastest parallel algorithm is fundamentally different than the fastest sequential one, and the parallel one performs more operations.
- We can speed up this program by having each processor compute the primes from  $1 \dots \sqrt{N}$ .
  - Why does doing extra computation make the code faster?

## Sieve of Eratosthenes

To find all primes  $\leq N$ :

- Let `MightBePrime = [2, 3, ..., N]`.
  - Let `KnownPrimes = []`.
  - while (`MightBePrime  $\neq []$` ) do
    - Loop invariant: `KnownPrimes` contains all primes less than the smallest element of `MightBePrime`, and `MightBePrime` is in ascending order. This ensures that the first element of `MightBePrime` is prime.
    - Let `P = first element of MightBePrime`.
    - Append `P` to `KnownPrimes`.
    - Delete all multiples of `P` from `MightBePrime`.
    - end
- See [http://en.wikipedia.org/wiki/Sieve\\_of\\_Eratosthenes](http://en.wikipedia.org/wiki/Sieve_of_Eratosthenes)

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 8 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 9 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 10 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 11 / 22

## Prime-Sieve in Erlang

## A More Efficient Sieve

- If  $N$  is composite, then it has at least one prime factor that is at most  $\sqrt{N}$ .
- This means that once we've found a prime that is  $\geq \sqrt{N}$ , all remaining elements of `Maybe` must be prime.
- Revised code:
 

```
% primes(N): return a list of all primes ≤ N.
primes(N) when is_integer(N) and (N < 2) -> [].
primes(N) when is_integer(N) -> lists:seq(2, N).
primes([]) -> [].
primes([P | Ps]) -> [P | primes([X || X <= N, X is not divisor(P)]).
```
- Main idea
  - Find primes from  $1 \dots \sqrt{N}$ .
  - Divide  $\sqrt{N} + 1 \dots N$  evenly between processors.
  - Have each processor find primes in its interval.
- We can speed up this program by having each processor compute the primes from  $1 \dots \sqrt{N}$ .
  - Why does doing extra computation make the code faster?

## Prime-Sieve: Parallel Version

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 12 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 13 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 14 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 15 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 16 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 17 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 18 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 19 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 20 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 21 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 22 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 23 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 24 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 25 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 26 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 27 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 28 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 29 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 30 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 31 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 32 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 33 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 34 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 35 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 36 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 37 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 38 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 39 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 40 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 41 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 42 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 43 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 44 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 45 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 46 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 47 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 48 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 49 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 50 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 51 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 52 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 53 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 54 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1, 2017 55 / 22

Mark Greenstreet Performance-Loss CS 418 – Feb. 1,

## Overhead: Summary

Overhead is loss of performance due to extra work that the parallel program does that is not performed by the sequential version. This includes:

- **Communication:** parallel processes need to exchange data. A sequential program only has one process; so it doesn't have this overhead.
- **Synchronization:** Parallel processes may need to synchronize to guarantee that some operations (e.g. file writes) are performed in a particular order. For a sequential program, this ordering is provided by the program itself.
- **Extra Computation:**
  - Sometimes it is more efficient to repeat a computation in several different processes to avoid communication overhead.
  - Sometimes the best parallel algorithm is a different algorithm than the sequential version and the parallel one performs more operations.
- **Extra Memory:** Data structures may be replicated in several different processes.

Mark Greenstreet      Performance Loss      CS 418 – Feb. 1, 2017      16 / 22

## Limited Parallelism

Sometimes, we can't keep all of the processors busy doing useful work.

- **Non-parallelizable code**  
The dependency graph for operations is narrow and deep.
- **Idle processors**  
There is work to do, but it hasn't been assigned to an idle processor.
- **Resource contention**  
Several processes need exclusive access to the same resource.

Mark Greenstreet      Performance Loss      CS 418 – Feb. 1, 2017      17 / 22

## Lecture Summary

### Causes of Performance Loss in Parallel Programs

- Overhead
  - Communication, [slide 5](#)
  - Synchronization, [slide 9](#)
  - Computation, [slide 10](#)
  - Extra Memory, [slide 15](#)
- Other sources of performance loss
  - Non-parallelizable code, [slide 18](#)
  - Idle Processors, [slide 19](#)
  - Resource Contention, [slide 20](#)

Mark Greenstreet      Performance Loss      CS 418 – Feb. 1, 2017      18 / 22

## Review Questions

- What is overhead? Give several examples of how a parallel program may need to do more work or use more memory than a sequential program.
- Do programs running on a shared-memory computer have communication overhead? Why or why not?
- Do message passing program have synchronization overhead? Why or why not?
- Why might a parallel program have idle processes even when there is work to be done?

Mark Greenstreet      Performance Loss      CS 418 – Feb. 1, 2017      19 / 22

## Idle Processors

- There is work to do, but processors are idle.
- Start-up and completion costs.
- Work imbalance.
- Communication delays.

## Models of Parallel Computation

Mark Greenstreet

CpSc 418 – Feb. 6, 2017

- The RAM Model of Sequential Computation
- Models of Parallel Computation
- An entertaining proof

 Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license  
<http://creativecommons.org/licenses/by/4.0/>

Mark Greenstreet      Performance Loss      CS 418 – Feb. 1, 2017      20 / 22

## Objectives

### The RAM Model

#### RAM = Random Access Machine

- Axioms of the model
  - Machines work on words of a "reasonable" size.
  - A machine can perform a "reasonable" operation on a word as a single step.
    - such operations include addition, subtraction, multiplication, division, comparisons, bitwise logical operations, bitwise shifts and rotates.
  - The machine has an unbounded amount of memory.
    - A memory address is a "word" as described above.
    - Reading or writing a word of memory can be done in a single step.

Mark Greenstreet      Performance Loss      CS 418 – Feb. 1, 2017      21 / 22

## The Relevance of the RAM Model

- If a single step of a RAM corresponds (to within a factor close to 1) to a single step of a real machine.
- Then algorithms that are efficient on a RAM will also be efficient on a real machine.
- Historically, this assumption has held up pretty well.
  - For example, `mergesort` and `quicksort` are better than `bubblesort` on a RAM and on real machines, and the RAM model predicts the advantage quite accurately.
  - Likewise, for many other algorithms
    - graph algorithms, matrix computations, dynamic programming, ...
    - hard on a RAM generally means hard on a real machine as well: NP complete problems, undecidable problems, ....
- The RAM model does not account for energy:
  - Energy is the critical factor in determining the performance of a computation.
  - The energy to perform an operation drops rapidly with the amount of time allowed to perform the operation.

Mark Greenstreet      Performance Loss      CS 418 – Feb. 1, 2017      22 / 22

## The CTA Model

#### CTA = Candidate Type Architecture

- Axioms of the model
  - A computer is composed of multiple processors.
  - Each processor has
    - Local memory that can be accessed in a single processor step (like the RAM model).
    - A small number of connections to a communications network.
  - There is a communication network connecting the processors.
    - the general model:
      - The communication network is a graph where all vertices (processors and switches) have bounded degree.
      - Each edge has an associated bandwidth and latency.
    - the simplified model:
      - Global actions have a cost of  $\lambda$  times the cost of local actions.
      - $\lambda$  is assumed to be "large".
      - the exact communication mechanism is not specified.

Mark Greenstreet      Performance Loss      CS 418 – Feb. 1, 2017      23 / 23

## The Irrelevance of the RAM Model

- The RAM model is based on assumptions that don't correspond to physical reality:
  - Memory access time is highly non-uniform.
  - Architects make heroic efforts to preserve the illusion of uniform access time fast memory –
    - caches, out-of-order execution, prefetching, ...
  - but the illusion is getting harder and harder to maintain.
    - Algorithms that randomly access large data sets run much slower than more localized algorithms.
    - Growing memory size and processor speeds means that memory hierarchy.
- The RAM model does not account for energy:
  - Energy is the critical factor in determining the performance of a computation.
  - The energy to perform an operation drops rapidly with the amount of time allowed to perform the operation.

Mark Greenstreet      Performance Loss      CS 418 – Feb. 1, 2017      24 / 24

## The (Ir)Relevance of the CTA Model

- Recognizing that communication is expensive is the one, most important point to grasp to understand parallel performance.
  - CTA highlights the central role of communication.
  - PRAM ignores it.
- The general model is parameterized by the communication network
  - Can we apply results from analysing a machine with a 3-D toroidal mesh to a machine with fat trees?
  - PRAM ignores it.
- The simple model neglects bandwidth issues
  - Messages are assumed to be "small".
  - But, bigger messages often lead to better performance.
  - If we talk about bandwidth, do we mean the bandwidth of each link?
  - Or, do we mean the bisection bandwidth?

Mark Greenstreet      Performance Loss      CS 418 – Feb. 1, 2017      25 / 25

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      2 / 23

## The LogP Model

- Learn about models of computation
  - Sequential: [Random Access Machine \(RAM\)](#)
  - Parallel
    - [Parallel Random Access Machine \(PRAM\)](#)
    - [Candidate Type Architecture \(CTA\)](#)
    - [Latency-Overhead-Bandwidth-Processors \(LogP\)](#)
- **An entertaining algorithm and its analysis**
  - If a model has invalid assumptions, then we can show that algorithm 1 is faster than algorithm 2, but in real life algorithm 2 is faster.
  - Valiant's algorithm also provides some mathematical entertainment.

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      3 / 23

## The Irrelevance of the PRAM Model

- The PRAM model is based on assumptions that don't correspond to physical reality:
  - Connecting  $N$  processors with memory requires a switching network.
    - Logic gates have bounded fan-in and fan-out.
    - $\Rightarrow$  any switch fabric with  $N$  inputs (and/or  $N$  outputs) must have depth of at least  $\log N$ .
    - This gives a lower bound on memory access time of  $\Omega(\log N)$ .
  - Processors exist in physical space
    - $N$  processors take up  $\Omega(N)$  volume.
    - The processor has a diameter of  $\Omega(N^{1/3})$ .
    - Switch fabric at a distance of at most  $c$  (the speed of light).
      - This gives a lower bound on memory access time of  $\Omega(N^{1/3})$ .

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      4 / 23

## The CTA Model

#### The CTA Model

- CTA = Candidate Type Architecture
- Axioms of the model
  - A computer is composed of multiple processors.
  - Each processor has
    - Local memory that can be accessed in a single processor step (like the RAM model).
    - A small number of connections to a communications network.
  - There is a communication network connecting the processors.
    - the general model:
      - The communication network is a graph where all vertices (processors and switches) have bounded degree.
      - Each edge has an associated bandwidth and latency.
    - the simplified model:
      - Global actions have a cost of  $\lambda$  times the cost of local actions.
      - $\lambda$  is assumed to be "large".
      - the exact communication mechanism is not specified.

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      5 / 23

## Why does g stand for "bandwidth"?

- **Marketing!**
  - What if we used **b** for "bandwidth"?
  - Need a catchy acronym with 't', 'o', 'b', and 'p' ...
    - got it: **BLOP**
    - but the marketing department vetoed it.
- **Why does g stand for "bandwidth"?**
  - What if we used **b** for "bandwidth"?
  - Need a catchy acronym with 't', 'o', 'b', and 'p' ...
    - got it: **BLOP**
    - but the marketing department vetoed it.

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      6 / 23

## Why does g stand for "bandwidth"?

- **Marketing!**
  - What if we used **b** for "bandwidth"?
  - Need a catchy acronym with 't', 'o', 'b', and 'p' ...
    - got it: **BLOP**
    - but the marketing department vetoed it.
- **Why does g stand for "bandwidth"?**
  - What if we used **b** for "bandwidth"?
  - Need a catchy acronym with 't', 'o', 'b', and 'p' ...
    - got it: **BLOP**
    - but the marketing department vetoed it.

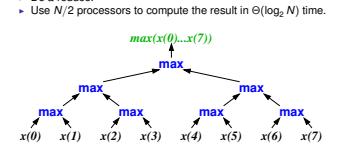
Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      7 / 23

## logP in practice

- The authors got some surprisingly good performance prediction for a few machines and a few algorithms by finding the "right" values for  $t$ ,  $o$ ,  $g$ , and  $P$  for each architecture.
- It's rare to get a model that comes to within 10-20% on several examples. So, this looked very promising.
- Since then, logP seems to be a model with more parameters than simplified CTA, but not particularly better accuracy.
- Good to know about, because if you meet an algorithms expert, they'll probably know that PRAM is unrealistic.
  - Then, you'll often hear "What about logP?" – the paper has lots of citations.
  - In practice, it's a slightly fancier way of saying "communication costs matter".

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      8 / 23

## Fun with the PRAM Model

- Finding the maximum element of an array of  $N$  elements.
  - The obvious approach
    - Do a reduce.
    - Use  $N/2$  processors to compute the result in  $\Theta(\log_2 N)$  time.
- 

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      9 / 23

## A Valiant Solution

### Valiant's algorithm, step 1

- Step 1:
  - Divide the  $N$  elements into  $N/3$  sets of size 3.
  - Assign 3 processors to each set, and perform all three pairwise comparisons in parallel.
  - Mark all the "losers" (requires a CRCW PRAM) and move the max of each set of three to a fixed location.
  - The PRAM operations in a bit more detail.
    - Initially, every element has a flag set to 1 that says "might be the max".
    - When  $\binom{k}{2}$  processors perform all of the pairwise comparisons of  $k$  values,
      - Each processor sets the flag for the smaller value to 0.
      - Note that several processors may write to the same location, but the CRCW allows this because they are all writing the same value.
    - One processor for each value checks if its flag is still set to 1.
      - The winner for the cluster is moved to a specific location.
      - The flag for that location is set to 1.
      - And now we're ready for subsequent rounds.

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      10 / 23

## Valiant's algorithm, step 1

- Step 1:
  - Divide the  $N$  elements into  $N/3$  sets of size 3.
  - Assign 3 processors to each set, and perform all three pairwise comparisons in parallel.
  - Mark all the "losers" (requires a CRCW PRAM) and move the max of each set of three to a fixed location.
  - The PRAM operations in a bit more detail.
    - Initially, every element has a flag set to 1 that says "might be the max".
    - When  $\binom{k}{2}$  processors perform all of the pairwise comparisons of  $k$  values,
      - Each processor sets the flag for the smaller value to 0.
      - Note that several processors may write to the same location, but the CRCW allows this because they are all writing the same value.
    - One processor for each value checks if its flag is still set to 1.
      - The winner for the cluster is moved to a specific location.
      - The flag for that location is set to 1.
      - And now we're ready for subsequent rounds.

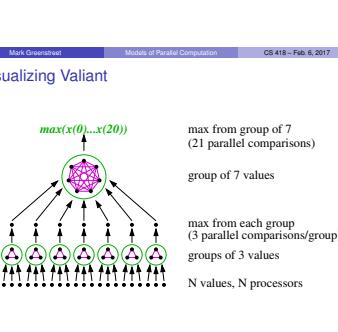
Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      11 / 23

## Valiant's algorithm, step 2

- We now have  $N/3$  elements left and still have  $N$  processors.
- We can make groups of 7 elements, and have 21 processors per group, which is enough to perform all  $\binom{7}{2} = 21$  pairwise comparisons in a single step.
- Thus, in  $O(1)$  time we move the max of each set to a fixed location. We now have  $N/21$  elements left to consider.
- The PRAM model is simple, and elegant, and many clever algorithms have been designed based on the PRAM model.
- It is also physically unrealistic:
  - As shown on [slide 7](#), logic gates have bounded fan-in and fan-out.
  - Implementing the processor-memory interconnect requires a logic network of depth  $\Omega(\log P)$ .
  - Therefore, access time must be  $\Omega(\log P)$ .
  - Each step of the PRAM takes  $\Omega(\log P)$  physical time.
- Valiant's  $O(\log \log N)$  algorithms takes  $O(\log \log \log N)$  physical time.
  - It's slower than doing a simple reduce.
  - And it uses lots of communication – think of all those  $\lambda$  penalties!
  - But it's very clever. ☺
- Valiant understood this and pointed these issues in his paper.
  - But there has still be extensive research on PRAM algorithms.
  - It's an elegant model, what can I say?

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      12 / 23

## Visualizing Valiant

- 
- max from group of 7 (21 parallel comparisons)
- group of 7 values
- max from each group (3 parallel comparisons/group)
- groups of 3 values
- N values, N processors

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      13 / 23

## Take-home message from Valiant's algorithm

- The PRAM model is simple, and elegant, and many clever algorithms have been designed based on the PRAM model.
- It is also physically unrealistic:
  - As shown on [slide 7](#), logic gates have bounded fan-in and fan-out.
  - Implementing the processor-memory interconnect requires a logic network of depth  $\Omega(\log P)$ .
  - Therefore, access time must be  $\Omega(\log P)$ .
  - Each step of the PRAM takes  $\Omega(\log P)$  physical time.
- Valiant's  $O(\log \log N)$  algorithms takes  $O(\log \log \log N)$  physical time.
  - It's slower than doing a simple reduce.
  - And it uses lots of communication – think of all those  $\lambda$  penalties!
  - But it's very clever. ☺
- Valiant understood this and pointed these issues in his paper.
  - But there has still be extensive research on PRAM algorithms.
  - It's an elegant model, what can I say?
- Simplified CTA reminds us that communication is expensive, but it doesn't explicitly charge for bandwidth.
- LogP accounts for bandwidth, but doesn't recognize that all bandwidth is not the same:
  - Communicating with an immediate neighbour is generally much cheaper than communicating with a distant machine.
  - Otherwise stated, the bisection bandwidth for real machines is generally much less than the per-machine bandwidth times the number of machines.
    - We can't have everyone talk at once at full bandwidth.
    - log<sub>2</sub> uses less bandwidth, but this is conservative, but it doesn't recognize the advantages of local communication.
- Both are based on a 10-20 year old machine model
  - That's ok, the papers are 18-25 years old.
  - Doesn't account for the heterogeneity of today's parallel computers:
    - multi-core on chip, faster communication between processors on the same board than across boards, etc.
- We'll use CTA because it's simple.
  - But recognize the limitations of any of these models.
- Getting a model of parallel computation that is all-purpose as the RAM is still a work-in-progress.

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      14 / 23

## Valiant's Algorithm, the remaining steps

### Valiant's Algorithm, run-time

- The sparsity is roughly squared at each step.
- It follows that the algorithm requires  $O(\log \log N)$ .
- Valiant showed a matching lower bound and extended the results to show merging is  $\Omega(\log \log N)$  and sorting is  $\Omega(\log N)$  on a CRCW PRAM.
- See [slide 27](#) for the details.

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      15 / 23

## Valiant's Algorithm, run-time

- The sparsity is roughly squared at each step.
- It follows that the algorithm requires  $O(\log \log N)$ .
- Valiant showed a matching lower bound and extended the results to show merging is  $\Omega(\log \log N)$  and sorting is  $\Omega(\log N)$  on a CRCW PRAM.
- See [slide 27](#) for the details.

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      16 / 23

## Summary

- Simplified CTA reminds us that communication is expensive, but it doesn't explicitly charge for bandwidth.
- LogP accounts for bandwidth, but doesn't recognize that all bandwidth is not the same:
  - Communicating with an immediate neighbour is generally much cheaper than communicating with a distant machine.
  - Otherwise stated, the bisection bandwidth for real machines is generally much less than the per-machine bandwidth times the number of machines.
    - We can't have everyone talk at once at full bandwidth.
    - log<sub>2</sub> uses less bandwidth, but this is conservative, but it doesn't recognize the advantages of local communication.
- Both are based on a 10-20 year old machine model
  - That's ok, the papers are 18-25 years old.
  - Doesn't account for the heterogeneity of today's parallel computers:
    - multi-core on chip, faster communication between processors on the same board than across boards, etc.
- We'll use CTA because it's simple.
  - But recognize the limitations of any of these models.
- Getting a model of parallel computation that is all-purpose as the RAM is still a work-in-progress.

Mark Greenstreet      Models of Parallel Computation      CS 418 – Feb. 6, 2017      17 / 23

**Preview**

- February 8: Parallel Sorting – The Zero-One Principle  
Reading: [https://en.wikipedia.org/w/index.php?title=Sorting\\_network&oldid=69011111](https://en.wikipedia.org/w/index.php?title=Sorting_network&oldid=69011111)
- February 10: Bitonic Sorting (part 1)  
Reading: [https://en.wikipedia.org/w/index.php?title=Bitonic\\_sorter&oldid=69011111](https://en.wikipedia.org/w/index.php?title=Bitonic_sorter&oldid=69011111)
- February 13: Family Day – no class
- February 15: Bitonic Sorting (part 2)  
Homework: HW 3 due (11:59pm), HW 4 goes out.
- February 17: Map-Reduce  
Homework: HW 3 due (11:59pm), HW 4 goes out.
- February 27: TBD
- March 1: Midterm
- March 3: GPU Overview  
Reading: [https://en.wikipedia.org/w/index.php?title=GPU\\_Computing\\_Era&oldid=69011111](https://en.wikipedia.org/w/index.php?title=GPU_Computing_Era&oldid=69011111)
- March 6: Intro to CUDA  
Reading: Kirk & Heav Ch. 2
- March 8: CUDA Threads, Part 1  
Reading: Kirk & Heav Ch. 3  
Homework: HW 4 earlybird (11:59pm)
- March 8: CUDA Threads, Part 2  
Homework: HW4 Due (11:59pm).

**Mark Greenstreet** Models of Parallel Computation CS 418 – Feb. 6, 2017 22 / 23

**Valiant Details**

round values remaining group size processors per group

|       |                     |                           |                                                        |
|-------|---------------------|---------------------------|--------------------------------------------------------|
| 1     | N                   | $2 \times 1 + 1 = 3$      | $3 = 3$ choose 2                                       |
| 2     | $\frac{N}{2}$       | $2 \times 3 + 1 = 7$      | $3 + 7 = 21 = 7$ choose 2                              |
| 3     | $\frac{N}{4}$       | $2 \times 21 + 1 = 43$    | $21 + 43 = 903 = \text{choose } 2$                     |
| 4     | $\frac{N}{8}$       | $2 \times 43 + 1 = 1,807$ | $903 + 1,807 = 1,801,721 = 1,807 \text{ choose } 2$    |
| ...   |                     |                           |                                                        |
| k     | $\frac{N}{2^k}$     | $2m_k + 1$                | $m_k(2m_k + 1) = (2m_{k+1} + 1) \text{ choose } 2$     |
| k + 1 | $\frac{N}{2^{k+1}}$ | $2m_{k+1} + 1$            | $m_k(2m_{k+1} + 1) = (2m_{k+1} + 1) \text{ choose } 2$ |

$m_k$  is the "sparsity" at round k:  
 $m_1 = 1$   
 $m_{k+1} = m_k(2m_k + 1) > 2m_k^2 > m_k^2$ .

Now note that  $m_{k+1} = m_k(2m_k + 1) > 2m_k^2 > m_k^2$ .  
Thus,  $\log(m_{k+1}) > 2\log(m_k)$ .  
For  $k \geq 3$ ,  $m_k > 2^{k-1}$ .  
Therefore, if  $N \geq 2$ ,  $k > \log(N) + 1 \Rightarrow m_k > N$ .

**Mark Greenstreet** Models of Parallel Computation CS 418 – Feb. 6, 2017 26 / 23

**Parallelizing Mergesort**  
We could use reduce?

**Mark Greenstreet** Models of Parallel Computation CS 418 – Feb. 6, 2017 27 / 23

**Parallelizing Mergesort**  
We could use reduce?

**Mark Greenstreet** Models of Parallel Computation CS 418 – Feb. 6, 2017 27 / 23

**Sorting Networks**

**Sorting Network for 2-elements**

```
in[1] --> a max --> out[1]
in[0] --> b min --> out[0]
```

**A Sorting Network for 3-elements**

```
in[2] --> a max --> out[2]
in[1] --> b min --> out[1]
in[0] --> b min --> out[0]
```

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 4 / 20

**Sorting Networks: Definition**  
Decision-tree version:

Let  $v$  be an arbitrary vertex of a decision tree, and let  $x_i$  and  $x_j$  be the variables compared at vertex  $v$ .  
A decision tree is a sorting network iff for every such vertex, the left subtree is the same as the right subtree with  $x_i$  and  $x_j$  exchanged.

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 8 / 20

**Compare-and-Swap Commutes with Monotonic Functions**

Compare-and-Swap commutes with monotonic functions.

Case  $x \leq y$ :

$$\begin{aligned} f(x) &\leq f(y), && \text{because } f \text{ is monotonic.} \\ \max(f(x), y) &= f(y), && \text{because } f(x) \leq f(y) \\ \max(f(x), y) &= f(\max(x, y)), && \text{because } x \leq y \end{aligned}$$

Case  $x \geq y$ : equivalent to the  $x \leq y$  case.

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 12 / 20

**Summary**

- Sequential sorting algorithms don't parallelize in an "obvious" way because they tend to have sequential bottlenecks.
  - Later, we'll see that we can combine ideas from sorting networks and sequential sorting algorithms to get practical, parallel sorting algorithms.
- Sorting networks are a restricted class of sorting algorithms
  - Based on compare-and-swap operations.
  - The parallelize well.
  - They don't have control-flow branches – this makes them attractive for architectures with large branch-penalties.
- The zero-one principle:
  - If a sorting-network sorts all inputs of 0s and 1s correctly, then it sorts all inputs correctly.
  - This allows many sorting networks to be proven correct by counting arguments.

**Review 1**

February 10: Bitonic Sorting (part 1)  
Reading: [https://en.wikipedia.org/w/index.php?title=Bitonic\\_sorter&oldid=69011111](https://en.wikipedia.org/w/index.php?title=Bitonic_sorter&oldid=69011111)

February 13: Family Day – no class

February 15: Bitonic Sorting (part 2)  
Homework: HW 3 due (11:59pm), HW 4 goes out.

February 17: Map-Reduce  
Homework: HW 3 due (11:59pm), HW 4 goes out.

February 27: TBD

March 1: Midterm

March 3: GPU Overview  
Reading: [https://en.wikipedia.org/w/index.php?title=GPU\\_Computing\\_Era&oldid=69011111](https://en.wikipedia.org/w/index.php?title=GPU_Computing_Era&oldid=69011111)

March 6: Intro to CUDA  
Reading: Kirk & Heav Ch. 2

March 8: CUDA Threads, Part 1  
Reading: Kirk & Heav Ch. 3  
Homework: HW 4 earlybird (11:59pm)

March 8: CUDA Threads, Part 2  
Homework: HW4 Due (11:59pm).

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 16 / 20

**For further reading**

- [Valiant1975] Leslie G. Valiant, "Parallelism in Comparison Problems," *SIAM Journal of Computing*, vol. 4, no. 3, pp. 348–355, (Sept. 1975).
- [Fortune1979] Steven Fortune and James Wyllie, "Parallelism in Random Access Machines," *Proceeding of the 11th ACM Symposium on Theory of Computing (STOC'79)*, pp. 114–118, May 1978.
- [Snyder1986] Lawrence Snyder, "Type architectures, shared memory, and the corollary of modest potential," *Annual review of computer science*, vol. 1, no. 1, pp. 289–317, 1986.
- [Guller1993] David Culler, Richard Karp, et al., "Log<sub>2</sub>-towards a realistic model of parallel computation," *ACM SIGPLAN Notices*, vol. 28, no. 7, pp. 1–12, (July 1993).

**EREW, CREW, and CRCW**

- EREW:** Exclusive-Read, Exclusive-Write
  - If two processors access the same location on the same step, then the machine fails.
- CREW:** Concurrent-Read, Exclusive-Write
  - Multiple machines can read the same location at the same time, and they all get the same value.
  - At most one machine can try to write a particular location on any given step.
  - If one processor writes to a memory location and another tries to read or write that location on the same step, then the machine fails.
- CRCW:** Concurrent-Read, Concurrent-Write
  - If two or more machines try to write the same memory word at the same time, then if they are all writing the same value, that value will be written. Otherwise (depending on the model),
    - the machine fails,
    - one of the writes "wins", or
    - an arbitrary value is written to that address.

**Mark Greenstreet** Models of Parallel Computation CS 418 – Feb. 6, 2017 29 / 29

**Mark Greenstreet** Models of Parallel Computation CS 418 – Feb. 6, 2017 29 / 29

**Mark Greenstreet** Models of Parallel Computation CS 418 – Feb. 6, 2017 24 / 25

**Mark Greenstreet** Models of Parallel Computation CS 418 – Feb. 6, 2017 25 / 25

**Sorting Networks**

**Mark Greenstreet** CS 418 – Feb. 8, 2017

**Parallelizing Mergesort**  
We could use reduce?

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 1 / 20

**Parallelizing Mergesort and/or quicksort**

- Sorting Networks
- The 0-1 Principle
- Summary

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license <http://creativecommons.org/licenses/by/4.0/>

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 2 / 20

**Parallelizing Mergesort**  
We could use reduce?

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 2 / 20

**Parallelizing Quicksort**  
How would you write a parallel version of quicksort?

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 3 / 20

**Parallelizing Quicksort**

**Mark Greenstreet** CS 418 – Feb. 8, 2017 3 / 20

**Sorting Networks: Definition**

**Structural version:**

- A sorting network is an acyclic network consisting of compare-and-swap modules.
  - Each primary input is connected either to the input of exactly one compare-and-swap module or to exactly one primary output.
  - Each compare-and-swap input is connected either to a primary input or to the output of exactly one compare-and-swap module.
  - Each compare-and-swap output is connected either to a primary output or to the input of exactly one compare-and-swap module.
  - Each primary output is connected either to the output of exactly one compare-and-swap module or to exactly one primary input.
- More formally, a sorting network is either
  - the identity network (no compare and swap modules),
  - a sorting network,  $S$  composed with a compare-and-swap module such that two outputs of  $S$  are the compare-to-the-and-swap, and the outputs of the compare-and-swap are the outputs of the new sorting network (along with the other outputs of the original network).

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 4 / 20

**Sorting Networks – Drawing**

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 5 / 20

**Sorting Networks – Examples**

Operations of the same color can be performed in parallel.

See: <http://pages.ripco.net/~jgambale/nw.html>

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 6 / 20

**The 0-1 Principle**

If a sorting network correctly sorts all inputs consisting only of 0s and 1s, then it correctly sorts inputs consisting of arbitrary comparable values.

- The 0-1 principle doesn't hold for arbitrary algorithms:
  - Consider the following linear-time "sort":
  - In linear time, count the number of zeros,  $n_0$ , in the array.
  - Set the first  $n_0$  elements of the array to zero.
  - Set the remaining elements to one.
  - This correctly sorts any array consisting only of 0s and 1s, but does not correctly sort other arrays.
- By restricting our attention to sorting networks, we can use the 0-1 principle.

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 7 / 20

**The 0-1 Principle: Proof Sketch**

We will show the contrapositive: if  $y$  is not sorted properly, then there exists an  $\bar{x}$  consisting of only 0s and 1s that is not sorted properly.

Choose  $i < j$  such that  $y_i > y_j$ .  
Let  $\bar{x}_i = 0$  if  $x_i < x_j$  and  $\bar{x}_i = 1$  otherwise.  
Clearly  $\bar{x}$  consists of only 0s and 1s.  
We will show that the sorting network does not sort correctly with input  $\bar{x}$ .

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 8 / 20

**Monotonicity Lemma**

Lemma: sorting networks commute with monotonic functions.

Let  $S$  be a sorting network with  $n$  inputs and  $n$  outputs.

- I'll write  $x_0, \dots, x_{n-1}$  to denote the inputs of  $S$ .
- I'll write  $y_0, \dots, y_{n-1}$  to denote the outputs of  $S$ .

Let  $f$  be a monotonic function.

- If  $x \leq y$ , then  $f(x) \leq f(y)$ .

The monotonicity lemma says

- applying  $S$  and then  $f$  produces the same result as
- applying  $f$  and then  $S$ .

Observation:  $f(X)$  when  $X < X_i \rightarrow 0$ ;  $f(X_i) \rightarrow 1$ .  $X$  is monotonic.

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 9 / 20

**The 0-1 Principle**

If a sorting network correctly sorts all inputs consisting only of 0s and 1s, then it correctly sorts inputs of any values.

I'll prove the contrapositive.

- If a sorting network does not correctly sort all inputs of 0s and 1s, then it does not correctly sort all inputs consisting only of 0s and 1s.
- Let  $S$  be a sorting network, let  $b$  be an input vector, and let  $y = S(b)$ , such that there exist  $i$  and  $j < i$  such that  $y_i > y_j$ .
- Let  $f(x) = 0, \quad \text{if } x < y_i \\ = 1, \quad \text{if } x \geq y_j$   
 $\bar{y} = S(f(x))$
- By the definition of  $f$ ,  $f(x)$  is an input consisting only of 0s and 1s.
- By the monotonicity lemma,  $\bar{y} = f(y) = 1 > 0 = f(y_i) = \bar{y}_i$
- Therefore,  $S$  does not correctly sort an input consisting only of 0s and 1s.

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 10 / 20

**Review 2**

Why don't traditional, sequential sorting algorithms parallelize well?

Try to parallelize another sequential sorting algorithm such as heap sort? What issues do you encounter?

Consider network sort-5(v2) from slide 6. Use the 0-1 principle to show that it sorts correctly?

- What are the inputs?
- What if one input has exactly one 1?
- What if the input has exactly two 1s?
- What if the input has exactly three 1s? Note, it may be simpler to think of this the input having exactly two 0s.
- What if the input has exactly four 1s? Five ones?

**Mark Greenstreet** CS 418 – Feb. 8, 2017 11 / 20

**Review 2**

sort-5 (v3) sort-5 (v4)

Consider the two sorting networks shown above. One sorts correctly; the other does not.

- Identify the network that sorts correctly, and prove it using the 0-1 principle.
- Show that the other network does not sort correctly by giving an input consisting of 0s and 1s that is not sorted correctly.

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 12 / 20

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 13 / 20

**Review 1**

February 10: Bitonic Sorting (part 1)  
Reading: [https://en.wikipedia.org/w/index.php?title=Bitonic\\_sorter&oldid=69011111](https://en.wikipedia.org/w/index.php?title=Bitonic_sorter&oldid=69011111)

February 13: Family Day – no class

February 15: Bitonic Sorting (part 2)  
Homework: HW 3 due (11:59pm), HW 4 goes out.

February 17: Map-Reduce  
Homework: HW 3 due (11:59pm), HW 4 goes out.

February 27: TBD

March 1: Midterm

March 3: GPU Overview  
Reading: [https://en.wikipedia.org/w/index.php?title=GPU\\_Computing\\_Era&oldid=69011111](https://en.wikipedia.org/w/index.php?title=GPU_Computing_Era&oldid=69011111)

March 6: Intro to CUDA  
Reading: Kirk & Heav Ch. 2

March 8: CUDA Threads, Part 1  
Reading: Kirk & Heav Ch. 3  
Homework: HW 4 earlybird (11:59pm)

March 8: CUDA Threads, Part 2  
Homework: HW4 Due (11:59pm).

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 14 / 20

**Review 1**

February 10: Bitonic Sorting (part 1)  
Reading: [https://en.wikipedia.org/w/index.php?title=Bitonic\\_sorter&oldid=69011111](https://en.wikipedia.org/w/index.php?title=Bitonic_sorter&oldid=69011111)

February 13: Family Day – no class

February 15: Bitonic Sorting (part 2)  
Homework: HW 3 due (11:59pm), HW 4 goes out.

February 17: Map-Reduce  
Homework: HW 3 due (11:59pm), HW 4 goes out.

February 27: TBD

March 1: Midterm

March 3: GPU Overview  
Reading: [https://en.wikipedia.org/w/index.php?title=GPU\\_Computing\\_Era&oldid=69011111](https://en.wikipedia.org/w/index.php?title=GPU_Computing_Era&oldid=69011111)

March 6: Intro to CUDA  
Reading: Kirk & Heav Ch. 2

March 8: CUDA Threads, Part 1  
Reading: Kirk & Heav Ch. 3  
Homework: HW 4 earlybird (11:59pm)

March 8: CUDA Threads, Part 2  
Homework: HW4 Due (11:59pm).

**Mark Greenstreet** Sorting Networks CS 418 – Feb. 8, 2017 15 / 20

## Review 3

I claimed that `max` and `min` can be computed without branches. We could work on the hardware design for a compare-and-swap module. Instead, I consider an algorithm that takes two "words" as arguments – each word is represented as a list of characters. The algorithm is supposed to output the two words, but in alphabetical order. For example:

```
See: http://www.ugrad.cs.ubc.ca/~cs418/2016/2/lecture/02-08/cas.erl
compareAndSwap(L1, L2) when isList(L1), isList(L2) ->
 compareAndSwap([L1], [L2]);
compareAndSwap([L1 | L2], [L3]) ->
 [lists:reverse(X), lists:reverse(X, L3)];
compareAndSwap([L1, L2], X) ->
 [lists:reverse(X), lists:reverse(X, L1)];
compareAndSwap([H1 | T1], [H2 | T2], X) when H1 == H2 ->
 compareAndSwap([T1, T2], [H1 | X]);
compareAndSwap([L1 | L2], [L3 | L4], X) when H1 < H2 ->
 [lists:reverse(X), lists:reverse(X, L3)];
compareAndSwap([L1, L2], X) ->
 [lists:reverse(X, L1), lists:reverse(X, L2)];
```

Show that `compareAndSwap` can be implemented as a scan operation.

Mark Greenstreet Sorting Networks CS 418 – Feb. 8, 2017 20 / 20

### Dividing the problem (part 1)

- For simplicity, assume each array has an even number of elements.
  - As we go on, we'll assume that each array has a power-of-two number of elements.
  - That's the easiest way to explain bitonic sort.
  - Note: the algorithm works for arbitrary array sizes.
    - See the [lecture slides from 2013](#).
- Divide each array in the middle?
  - If  $A$  has  $N$  elements and  $A_i$  are ones.
  - How many ones are in  $A[0..(N/2)-1]$ ?
  - How many ones are in  $A[(N/2)..N-1]$ ?
- Taking every other element?
  - How many ones are in the  $A[0, 2, \dots, N-2]$ ?
  - How many ones are in the  $A[1, 3, \dots, N-1]$ ?
- Other schemes?

Mark Greenstreet Bitonic Sort CS 418 – Feb. 10, 2017 4 / ??

### Bitonic Merge – big picture

- Bitonic merge produces a monotonic sequence from an bitonic input.
  - Given two sorted sequences,  $A$  and  $B$ , note that
- $$X = A + \text{reverse}(B)$$
- is bitonic.
- We don't require the lengths of  $A$  or  $B$  to be powers of two.
  - If fact, we don't even require that  $A$  and  $B$  have the same length.
  - Divide  $X$  into  $X_0$  and  $X_1$ , the even-indexed and odd-indexed subsequences.
    - $X_0$  and  $X_1$  are both bitonic.
    - The number of 1s in  $X_0$  and  $X_1$  differ by at most 1.
  - Use bitonic merge (recursion) to sort  $X_0$  and  $X_1$  into ascending order to get  $Y_0$  and  $Y_1$ .
    - HowManyOnes( $Y_0$ ) = HowManyOnes( $X_0$ ), and HowManyOnes( $Y_1$ ) = HowManyOnes( $X_1$ ).
    - Therefore, the number of 1s in  $Y_0$  and  $Y_1$  differ by at most 1.
    - This is an "easy" case from [slide 3](#).

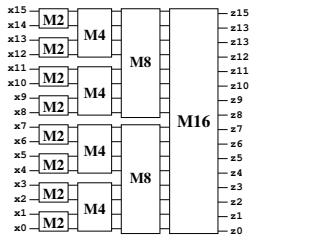
Mark Greenstreet Bitonic Sort CS 418 – Feb. 10, 2017 8 / ??

### The complexity of bitonic merge

- We'll count the compare-and-swap operations
  - Is it OK to ignore reversing one array, concatenating the arrays, separating the even- and odd-indexed elements, and recombining them?
  - Yes. The number of these operations is proportional to the number of compare-and-swaps.
  - Yes. Even better, in the next lecture, we'll show how to eliminate most of these data-shuffling operations.
- A bitonic merge of  $N$  elements requires:
  - two bitonic merges of  $\lfloor N/2 \rfloor$  items (if  $N > 2$ )
  - $\lfloor N/2 \rfloor$  compare-and-swap operations.
- The total number of compare and swap operations is  $O(N \log N)$ .

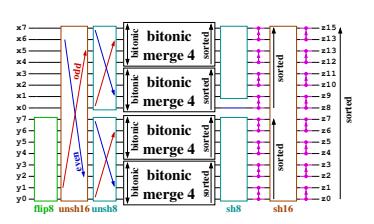
Mark Greenstreet Bitonic Sort CS 418 – Feb. 10, 2017 12 / ??

### Parallelizing Mergesort



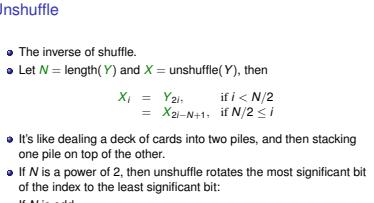
Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 2 / 17

### Bitonic Merge



Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 3 / 17

### Bitonic Merge



Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 3 / 17

### Unshuffle

- The inverse of shuffle.
  - Let  $N = \text{length}(Y)$  and  $X = \text{unshuffle}(Y)$ , then
- $$X_i = Y_{2^j}, \quad \text{if } i < N/2 \\ = X_{2^j-N+1}, \quad \text{if } N/2 \leq i$$
- It's like dealing a deck of cards into two piles, and then stacking one pile on top of the other.
- If  $N$  is a power of 2, then unshuffle rotates the most significant bit of the index to the least significant bit:
  - If  $N$  is odd,

Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 3 / 17

### Bitonic Sort

Mark Greenstreet

CpSc 418 – Feb. 10, 2017

- Merging
- Shuffle and Unshuffle
- The Bitonic Sort Algorithm
- Summary

- I know that some of the links in the electronic version are broken. I know that it would be nice if I complete the final slides. I will post to piazza when this is done.

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license.

<http://creativecommons.org/licenses/by/4.0/>

### Parallelizing Mergesort

Mark Greenstreet

CpSc 418 – Feb. 10, 2017

- We looked at this in the [Feb. 8 lecture](#).
- The challenge is the merge step:
  - Can we make a parallel merge?

Mark Greenstreet Bitonic Sort CS 418 – Feb. 10, 2017 1 / ??

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

CpSc 418 – Feb. 10, 2017

Mark Greenstreet

Bitonic Sort

</div

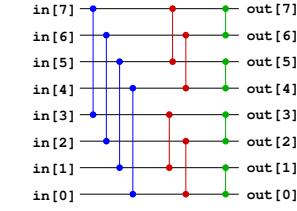
## The first compare-and-swap

- The first compare-and-swap operates on  $y(b_3, b_2, b_1, 0)$  and  $y(b_3, b_2, b_1, 1)$ , for all 8 choices of  $b_3$ ,  $b_2$ , and  $b_1$ .
- This corresponds to a compare-and-swap of  $x(b_0, b_1, b_2, 0)$  with  $x(b_0, b_1, b_2, 1)$ .
  - I'll call the result of the compare-and-swap  $z$  where
    - $z(b_3, b_2, b_1, 0) = \min(y(b_3, b_2, b_1, 0), y(b_3, b_2, b_1, 1))$
    - $z(b_3, b_2, b_1, 1) = \max(y(b_3, b_2, b_1, 0), y(b_3, b_2, b_1, 1))$
  - And I'll write  $\tilde{z}$  for  $z$  with "x" indexing:
    - $\tilde{z}(b_3, b_2, b_1, b_0) = z(b_3, b_2, b_1, b_0)$
    - $\tilde{z}(b_3, b_2, b_1, 0) = \min(x(b_3, b_2, b_1, 0), x(b_3, b_2, b_1, 1))$
    - $\tilde{z}(b_3, b_2, b_1, 1) = \max(x(b_3, b_2, b_1, 0), x(b_3, b_2, b_1, 1))$
  - These are comparisons with a "stride" of 8 for  $(x)$ .

## The first shuffle

- The first shuffle takes  $z$  as an input and I'll call the output  $w$ .
- The first shuffle is a **shuffle-4**; so
  - $w[i] = z(i\text{rot}_1(1, 2))$
  - Equivalently,  $w(b_3, b_2, b_1, b_0) = z(b_3, b_2, b_0, b_1)$ .
- Let
 
$$\begin{aligned} \tilde{w}(b_3, b_2, b_1, b_0) &= w(b_3, b_1, b_2, b_0) \\ &= z(b_3, b_2, b_3, b_0) \\ &= \tilde{z}(b_3, b_2, b_1, b_0) \end{aligned}$$
- The second stage of compare-and-swap modules operates on
  - $w(b_3, b_2, b_1, 0)$  and  $w(b_3, b_2, b_1, 1)$
  - Equivalently,  $\tilde{w}(b_3, b_2, b_1, 0)$  and  $\tilde{w}(b_3, b_2, b_1, b_0)$ .
  - These are comparisons with a stride of 4 for  $\tilde{z}$  and  $\tilde{w}$ .

## The "Textbook" Diagram



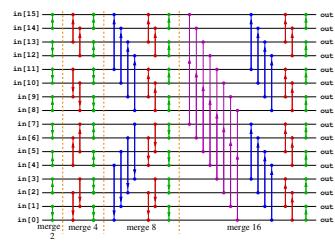
## Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 9 / 17

### Flipping Out

What should we do about the flips?

- Push them back (right-to-left) through the network
  - Keep track of how many flips we've accumulated.
  - Sort up for an even number of flips.
  - Sort down for an odd number of flips.
- Flip the wiring in the bottom half of each unshuffle.
- In practice:
  - Do the one that's easier for your implementation.

## Bitonic Sort



## Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 10 / 17

### Related Algorithms

- Counting Networks
  - How to match servers to requests.
- FFT
  - The Platonic Ideal of a Divide-and-Conquer Algorithm
  - Used for speech processing, signal processing, and lots of scientific computing tasks.

## Map-Reduce

### Mark Greenstreet & Ian M. Mitchell CPSC 418 – February 27, 2017

- Problem Domain: Large-Scale Data Analysis**
- The Map-Reduce Pattern**
- Implementation Issues**
- Results**

Unsplash image or clip, these slides are copyright 2017 by Mark Greenstreet & Ian Mitchell and are made available under the terms of the Creative Commons Attribution 4.0 International license  
http://creativecommons.org/licenses/by/4.0/

## The MapReduce Pattern

Slight generalization of description from [Dean & Ghemawat 2008].

- All data is represented as collections of **(Key, Value)** pairs.
- map**
  - For each **(Key1, Value1)** pair of the input, user code produces a collection of **(Key2, Value2)** pairs for the output.
- shuffle**
  - All **(Key2, Value2)** pairs with the same **Key2** are combined into a **[Key2, list of Value2]** result and sorted by **Key2**.
- reduce**
  - For each **(Key2, list of Value2)**, user code produces a **(Key2, Value3)** result (where **Value3** might be a list itself).

Apache Hadoop is an open-source implementation of this basic framework.

## Greenstreet & Mitchell Map-Reduce CPSC 418 – Feb. 27, 2017 4 / 21

- In Erlang, **wtree:reduce()** is designed to spread the computation of a reduction across many workers.
  - Implementation maximizes parallelism for a single reduce operation.
  - Collection and combination of data occurs in **combine()** functions.
  - Note that the **leaf()** function can perform a map operation before the reduction, so no loss of generality.
- In MapReduce, many independent reductions (one for each **Key2**) are spread across many workers, but each reduction is performed by a single worker.
  - Implementation emphasizes fault tolerance and disk-based data storage.
  - Collection (but not combination) of data occurs in shuffle step.
  - If reduction is too big for a single worker, user must change the intermediate **(Key2, Value2)** representation and/or break the problem into multiple MapReduce steps.

## MapReduce: Word-Count

Example from [Dean & Ghemawat 2008] revised.

- Input data:**
  - Key1** is the document name.
  - Value1** is document text.
- map:**
  - Key2** is a word.
  - Value2** is the count of times it appears in the document.
- shuffle:**
  - Collects counts from all documents for each word (**Word**, **list of Counts**).
- reduce:**
  - Value3** is the sum of all counts; in other words, the total number times the word appears in all documents.

The user then invokes the **MapReduce** function.

## Greenstreet & Mitchell Map-Reduce CPSC 418 – Feb. 27, 2017 8 / 21

- Map workers know the number of reduce workers  $R$ .
  - Each **(Key2, Value2)** is written to a different file according to **hash(Key2) mod R**.
  - The master tells the reduce worker which file to read from each map worker.
- Later versions of MapReduce utilized more complex or even user-specific hashing; for example to:
- Better balance size of reduce problems.
  - Reduce network traffic and/or simplify sorting during shuffle step.
  - Cluster certain **Key2** tuples onto the same reduce workers.

## Fault Tolerance

Bad things happen: Failed disks, partitioned networks, power shortages, ...

- Key Idea:**
  - The **map** and **reduce** operations are based on functional programming ideas: referential transparency and no side-effects.
  - If a worker crashes, it is as if it never existed.
  - The master can restart the task on another machine.
- The master periodically pings the tasks, and restarts dead ones.
- Map tasks produce only temporary files, so if a completed map task fails before informing the master then it must be re-executed.
- Reduce tasks produce files in the distributed file system (redundant and fault-tolerant), so no need to re-execute.

## Performance Issues

- The master attempts to schedule map tasks on the processor whose local disk holds the split being processed, or are nearby (by the network connections).
- Shuffle moves data from many map tasks to many reduce tasks. Easily saturates the cross-section bandwidth of the network.
- For good performance, the map tasks should be filters that output much less data than they read.
  - Often not true of the "natural" intermediate representation (such as curriculum problem above).
  - Fewer distinct **Key2** values means less parallelism in reduce tasks.
- Can often reduce intermediate data size by partial reduction in the map workers.
  - May change the semantics of MapReduce (but not if reduce is associative and commutative).

## MapReduce is Designed for BIG Data

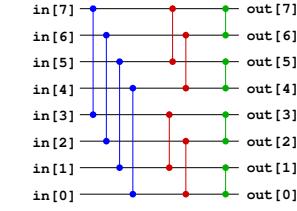
- Communication between standard linux machines with generic networks takes milliseconds.
- Reading large disk files takes seconds.
- The task needs to be big enough to justify these overheads:
  - Equivalent to increasing  $\lambda$  by a few orders of magnitude.
  - MapReduce makes sense if the task is disk-limited and harnessing a few thousand disks provides the necessary disk bandwidth.
    - Think of it as "disk parallelism" instead of "CPU parallelism".
    - Note: big-data companies like Amazon, Facebook and Google are moving to using FLASH memory and DRAM instead of disks, exactly because of these I/O bottlenecks.
  - Or if you have a really huge data set the compute time may dominate all of these overheads.

## The rest of the merge

- In the same way, the third stage of compare-and-swap modules operates with a stride of 2 for  $i$  indices.
- And the final stage has a stride of 1.
- More generally, to merge two sequences of length  $2^L$ .
  - Flip the lower sequence.
  - Or, just sort it in reverse in the first place.
  - Perform compare-and-swap operations with stride  $L$ .
  - Perform compare-and-swap operations with stride  $L/2$  – note that these operate on pairs of elements whose index bits differ in the  $L/2$  bit, and all of their other index bits are the same.
  - Perform compare-and-swap operations with stride  $L/4$ , ...
  - Perform compare-and-swap operations with stride 1 – this compares the element at  $2+i$  with the element at  $2+i+1$  for  $0 \leq i < 2^L$ .
- Done!

## Practical performance

## The "Textbook" Diagram



## Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 11 / 17

### Bitonic Sort in practice

- Sorting networks can be used to design practical sorting algorithms.
- To sort  $N$  values with  $P$  processors:
  - Divide input into  $2P$  segments of length  $\frac{N}{2}$ .
  - Each processor sorts its pair of segments into one long segment.
  - The sorted segments are the inputs to the sorting network.
  - Now, follow the actions of the sorting network:
    - Processor 1 handles row 2 and  $2+1$  of the sorting network.
    - Each compare-and-swap is replaced with "merge two sorted sequences and split into top half and bottom half".
    - When the sorting network has a compare-and-swap between rows  $2i$  and  $2i+1$ , each processor handles it locally.
    - Processor 1 sorting network has a compare-and-swap between rows  $2i$  and  $2i+1$ . This means processor 1 sends the upper half of its data to processor  $K-1$  (processor 1), and processor  $1-K$  sends the lower half of its data to processor 1. Both perform merges.
    - Note if the compare-and-swap was flipped, then flip "upper-half" and "lower half".

## Practical performance

### Complexity

- Total number of comparisons:  $O(N \log N \log^2 P)$ .
- Time:  $O(\frac{N}{P} (\log N + \log^2 P))$ , assuming each processor sorts  $N/P$  elements in  $O(N/P \log(N/P))$  time and merges two sequences of  $N/P$  elements in  $O(N/P)$  time.

### Remarks:

- The idea of replacing compare-and-swap modules with processors that can perform merge using an algorithm optimized for the processor, is an extremely powerful and general idea. It is used in the design of many practical parallel sorting algorithms.
- Sorting networks are becoming the SIMD branches:
  - Ideal for SIMD machines that can't really branch.
  - Need to experiment some to see the trade-offs of branch-divergence vs. higher asymptotic complexity on a GPU.

## Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 12 / 17

## Problem Domain: Large-Scale Data Analysis

### Data analysis requires:

- Fetching the relevant records.
- Performing analysis of related records.
- Summarizing the results.

### Example: word frequency in documents

- How do 200-level courses impact success in 400-level courses?
- Look at all transcripts.
- Analyze relationships for (2XX, 4YY) pairs.

- Google noted that each such problem was getting its own custom code.
  - All of that code development is expensive.
  - Often required similar rewrites when underlying system changed.

## Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 13 / 17

### Portrait of a Data Centre

### MapReduce

- Sketch of a big data center:
  - Tens of thousands of machines, each with its own disk(s).
  - Distributed file system – what is that?
  - Commodity networks and routers.
    - Each machine has a network interface (e.g. 10Gb ethernet)
    - Cross-section bandwidth is **way** smaller than the number of machines times the per-machine bandwidth.
- Scale: Is it big that there will be **failures**: chips, cores, memory, disks, network interfaces, switches, ...
  - If the average lifetime of a machine is five years, then 10,000 machine data center will have a failure every four hours.
  - Even without failure, maintenance will take machines offline.
- Need to analyse the huge data sets stored on these machines.

## Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 14 / 21

### MapReduce: Curriculum

- Input data (Key1, Value1):**
  - Key1** is the student number for the transcript
  - Value1** is a list of **(CourseNumber, Grade)** pairs.
- map:** For each 200-level course and for each 400-level course in the transcript:
  - Key2** is the pair **(Course200, Course400)**.
  - Value2** is the pair **(Grade200, Grade400)**.
- shuffle:**
  - Collect matching course pairs from all students **(Course200, Course400)**, **(Grade200, Grade400)**, **(Grade200, Grade400)**, ...).
- reduce:**
  - Value3** is (for example) the sample Pearson correlation coefficient  $r$  for the data set **(Grade200, Grade400)**, **(Grade200, Grade400)**, ...).
    - More complex analyses could be performed.
- Result is a bunch of temporary files holding **(Key2, Value2)** pairs.

## Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 15 / 21

### Execution (Part 1)

- The **MapReduce** function spawns  $M$  map worker,  $R$  reduce workers, and one master.
- Each map worker:
  - Is assigned fragment(s) of the input file by the master – these fragments are called **splits**.
  - Reads a **(Key1, Value1)** record from an assigned split.
  - Runs user **map** code on that record.
  - Writes the **(Key2, Value2)** result to a temporary file.
  - Repeats until all records in the split are completed.
  - Repeats until all assigned splits are completed.
  - Notifies the master when it is done.
- Result is a bunch of temporary files holding **(Key2, Value2)** pairs.

## Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 16 / 21

### Execution (Part 2)

- Start from temporary files holding **(Key2, Value2)** pairs.
- Shuffle:**
  - Each reduce worker is assigned **Key2** value(s).
  - Corresponding **Value2** lists are taken from map worker's temporary output files and written to reduce worker's temporary input files.
  - Reduce worker receives **(Key2, list of Value2)** pairs sorted by **Key2**.
- Each reduce worker:
  - Reads a **(Key2, list of Value2)** record from a temporary file.
  - Runs user **reduce** code on that record.
  - Writes the **(Key2, Value3)** result to a file.
  - Notifies the master when it is done.
- When all the reduce computations are complete, the master sends a message to wake up the user process, and the **MapReduce** function returns.

## Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 17 / 21

### Work Stealing

- Sometimes a worker is slow – **stragglers**.
- If the **MapReduce** task is near completion, the master assigns straggler tasks to idle processors.
- These are called **backup tasks**.
- Either the original or the backup process can complete the task.
- In practice, this **work stealing** by backup tasks:
  - Only adds a few percent to the total compute resources used.
  - Can result in substantial performance improvements: The paper reported a 44% slow-down when the sort benchmark was run without backup tasks.

## Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 18 / 21

### Results (Part 1)

- Achieves impressive performance on massive data sets 2008–2013?
- Report in [Dean & Ghemawat, 2008]: Good performance on ~ 2000 machines: **grep** and **sort** work through  $10^{10}$  100-byte records (1TB) in minutes.
- Google estimates ~ 20PB / day in total MapReduce processing in January 2008.
- Google research blog reports sorting  $10^{13}$  100-byte records (1PB) on ~ 4000 machines (and ~ 48,000 disks) in six hours in November 2008, then 33 minutes for 1 PB or 6 hours for 10 PB on 8000 machines in September 2011.
- Open source implementation in Hadoop widely available as a cloud service.
- Many example algorithms documented; for example, search for "map reduce" on <http://scholar.google.ca>.

## Mark Greenstreet Bitonic Sort CS 418 – Feb. 15, 2017 19 / 21

### Results (Part 2)

- Big data processors are now moving away from MapReduce.
- Framework emphasizes batch processing of data residing in distributed file system, which limits flexibility and efficiency.
- "We don't really use MapReduce anymore" (Urs Hölzle, senior vice president of technical infrastructure at Google speaking at Google I/O conference in 2014)
- Machine learning project Apache Mahout is shifting away from MapReduce algorithms to alternatives such as Apache Spark.
  - Worth noting: Spark also uses a functional programming model with referential transparency.

## Summary

- MapReduce is a parallel programming pattern
  - Data are represented by collections of `(Key, Value)` pairs.
  - User provides `map` to take input `(Key1, Value1)` pairs to intermediate `(Key2, Value2)` pairs.
  - Shuffle step collects intermediate data into `(Key2, [list of Value2])` pairs and sorts it by `Key2`.
  - User provides `reduce` to take sorted `(Key2, [list of Value2])` pairs into `(Key2, Value3)` pairs.
- Details of the parallel implementation are handled by the MapReduce API:
  - Creating workers processes, delivering input and output files, shuffling intermediate data between map and reduce workers.
  - Handling failures and slow nodes.
- Performance is often bandwidth limited.
  - Locality matters: perform `map` on the machine with the data.
  - If `map` is an effective filter (bytes out < bytes in), then we can reduce the impact of network congestion.
  - In practice, user must choose `(Key2, Value2)` representation wisely to trade-off shuffle effort for reduce parallelism.

## Review: Properties of MapReduce

## Review: Example of a MapReduce Problem

## Introduction to GPGPUs

Greenstreet & Mitchell Map Reduce CS 418 – Feb. 27, 2017 20 / 21

Greenstreet & Mitchell Map Reduce CS 418 – Feb. 27, 2017 21 / 21

Greenstreet & Mitchell Map Reduce CS 418 – Feb. 27, 2017 22 / 21

Mark Greenstreet CpSc 418 – Mar. 3, 2017

### Before the first GPU

### 1989: The SGI Geometry Engine

### 1989: The SGI Geometry Engine

### Why is dedicated hardware so much faster?

Early 1980's: bit-blt hardware for simple 2D graphics.
 

- Draw lines, simple curves, and text.
- Fill rectangles and triangles.
- Color used a "color map" to save memory:
  - bit-wise logical operations on color map indices!

- Basic rendering: coordinate transformation.
  - Represent a 3D point with a 4-element vector.
  - The fourth element is 1, and allows translations.
  - Multiply vector by matrix to perform coordinate transformation.
- Dedicated hardware is **much more efficient** than a general purpose CPU for matrix-vector multiplication.
  - For example, a 32 × 32 multiplier can be built with  $32^2 = 1024$  one-bit multiplier cells.
    - A one-bit multiplier cell is about 50 transistors.
    - That's about 50K transistors for a very simple design.
    - 30K is quite feasible using better architectures.

- Basic rendering: coordinate transformation.
- Dedicated hardware is **much more efficient** than a general purpose CPU for matrix-vector multiplication.
  - For example, a 32 × 32 multiplier can be built with  $32^2 = 1024$  one-bit multiplier cells.
    - A one-bit multiplier cell is about 50 transistors.
    - That's about 50K transistors for a very simple design.
    - 30K is quite feasible using better architectures.
- The 80486DX was also built in 1989.
  - The 80486DX was a 1.2MHz, 16-bit, 16MHz, 13MPflops.
  - That's equal to 24 dedicated multipliers.
  - 16 multiply-and-accumulate units running at 50MHz (easy in the same 1<sub>μ</sub>-process) produce 1.6GFlops!

Mark Greenstreet

Introduction to GPGPUs

CS 418 – Mar. 3, 2017 1 / 25

### Building a better multiplier

### Why is dedicated hardware so much faster?

### The fundamental challenge of graphics

Mark Greenstreet

Introduction to GPGPUs

CS 418 – Mar. 3, 2017 4 / 25

- Simple multiplier has latency and period of  $2N$ .
- Pipelining: add registers between rows.
  - The period is  $N$ , but the latency is  $N^2$ .
  - The bottleneck is the time for carries in each row.
- Use carry-lookahead adders (compute carries with a scan)
  - period is  $\log N$ , the latency is  $N/\log N$ .
  - but the hardware is more complicated.
- Use carry-save adders and one carry-lookahead at the end
  - each adder in the multiplier forwards its carries to the next adder.
  - the final adder resolves the carries.
  - period is latency is  $N$ .
  - and the hardware is **way simpler** than a carry-lookahead design
- Graphics and many scientific and machine learning computations are very tolerant of latency.

- Example: matrix-vector multiplication
- addition and multiplication are "easy".
  - it's the rest of the CPU that's complicated and the usual performance bottleneck
    - memory read and write
    - instruction fetch, decode, and scheduling
    - pipeline control
    - handling exceptions, hazards, and speculation
    - etc.
  - GPU architectures amortize all of this overhead over a lot of execution units.

- Human vision isn't getting any better.
- Once you can perform a graphics task at the limits of human perception (or the limits of consumer budget for monitors), then there's no point in doing it any better.
  - Rapid advances in chip technology meant that coordinate transformations (the specialty of the SGI Geometry Engine) were soon as fast as anyone needed.
  - Graphics processors have evolved to include more functions. For example,
    - Shading
    - Texture mapping
  - This led to a change from hardwired architectures, to programmable ones.

Mark Greenstreet

Introduction to GPGPUs

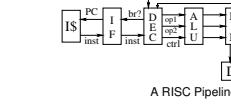
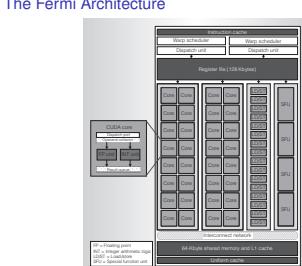
CS 418 – Mar. 3, 2017 3 / 25

### The Fermi Architecture

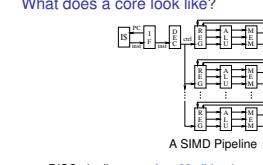
### What does a core look like?

### What does a core look like?

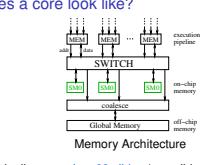
### What does a core look like?



- RISC pipeline: see Jan. 23 slides (e.g. slides 5ff.)
  - Instruction fetch, decode and other control takes much more power than actually performing ALU and other operations!
- SIMD: Single-Instruction, Multiple-Data
  - Each pipeline has its own registers and operates on separate data values.
  - Commonly, pipelines access **adjacent** memory locations.
  - Great for operating on matrices, vectors, and other arrays.
- What about memory?



- RISC pipeline: see Jan. 23 slides (e.g. slides 5ff.)
- SIMD: Single-Instruction, Multiple-Data
  - What about memory?
- Each pipeline switched between cores: see Jan. 23 slides (e.g., slide 3)
- Off-chip references are "coalesced": the hardware detects reads from (or writes to) consecutive locations and combines them into larger, block transfers.



### More about GPU Cores

### Lecture Outline

### Execution Model: Functions

### Execution Model: Memory

- Execution pipeline can be very deep – 20-30 stages.
  - Many operations are floating point and take multiple cycles.
  - A floating point unit that is deeply pipelined is easier to design, can provide higher throughput, and use less power than a lower latency design.
- No bypasses
  - Instructions block until instructions that they depend on have completed execution.
  - GPUs rely on extensive multi-threading to get performance.
- Banches use **predicated execution**:
  - Execute the then-branch code, disabling the "else-branch" threads.
  - Execute the else-branch code, disabling the "then-branch" threads.
  - The order of the two branches is unspecified.
- Why?
  - All of these choices optimize the hardware for graphics applications.
  - To get good performance, the programmer needs to understand how the GPGPU executes programs.

- GPUs
  - been there, done that.
- CUDA – we are here!
  - Execution Model
  - Memory Model
  - Code Snippets

- A CUDA program consists of three kinds of functions:
  - Host functions:
    - callable from code running on the host, but not the GPU.
    - run on the host CPU.
    - In CUDA C, these look like normal functions.
  - Device functions:
    - callable from code running on the GPU, but not the host.
    - run on the GPU.
    - In CUDA C, these are declared with a `__device__` qualifier.
  - Global functions:
    - called by code running on the host CPU,
    - they execute on the GPU.
    - In CUDA C, these are declared with a `__global__` qualifier.

- Host memory: DRAM and the CPU's caches
  - Accessible to host CPU but not to GPU.
- Device memory: GDDR DRAM on the graphics card.
  - Accessible by GPU.
  - The host can initiate transfers between host memory and device memory.
- The CUDA library includes functions to:
  - Allocate and free device memory.
  - Copy blocks between host and device memory.
  - BUT** host code can't read or write the device memory directly.

### Structure of a simple CUDA program

### saxpy: device code

### saxpy: host code (part 1 of 5)

### saxpy: host code (part 2 of 5)

- A `__global__` function to called by the host program to execute on the GPU.
  - There may be one or more `__device__` functions as well.
- One or more host functions, including `main` to run on the host CPU.
  - Allocate device memory.
  - Copy data from host memory to device memory.
  - "Launch" the device kernel by calling the `__global__` function.
  - Copy the result from device memory to host memory.
- We'll do the `saxpy` example from the paper.
  - `saxpy` = "Scalar a times x plus y".

```
int main(argc, char **argv) {
 float n = atoi(argv[1]);
 float *x, *y;
 float a;
 a = 2.5;
 n = 256;
 x = (float *)malloc(n * sizeof(float));
 y = (float *)malloc(n * sizeof(float));
 for(int i = 0; i < n; i++) {
 x[i] = i;
 y[i] = i;
 }
 __global__ void saxpy(uint n, float a, float *x, float *y) {
 uint i = blockIdx.x * blockDim.x + threadIdx.x;
 if(i < n) {
 y[i] = a*x[i] + y[i];
 }
 }
}
```

```
int main(void) {
 ...
 float a;
 a = 2.5;
 n = 256;
 x = (float *)malloc(n * sizeof(float));
 y = (float *)malloc(n * sizeof(float));
 for(int i = 0; i < n; i++) {
 x[i] = i;
 y[i] = i;
 }
 ...
}
```

- Declare variables for the arrays on the host and device.
- Allocate and initialize values in the host array.

### saxpy: host code (part 3 of 5)

### saxpy: host code (part 4 of 5)

### saxpy: host code (part 5 of 5)

### Threads and blocks

- Invoke the code on the GPU:
  - `cudaMemcpy(x, dev_x, size, cudaMemcpyDeviceToHost);` says to create  $[n/256]$  blocks of threads.
  - Each block consists of 256 threads.
  - See slide 22 for an explanation of threads and blocks.
  - The pointers to the arrays (in device memory) and the values of `n` and `a` are passed to the threads.
  - Copy the result back to the host.

```
for(int i = 0; i < n; i++) { // check the result
 if(y[i] != a*x[i] + y[i]) {
 fprintf(stderr, "ERROR: i=%d, a[%d]=%f, b[%d]=%f, c[%d]=%f\n",
 i, a[i], b[i], c[i]);
 exit(-1);
 }
}
printf("The results match!\n");
}
```

- Check the results.
- Clean up.
- We're done.

- Our example created  $\lceil \frac{n}{256} \rceil$  blocks with 256 threads each.
- The GPU hardware has a pool of running threads.
  - Each thread has a "next instruction" pending execution.
  - If the next instruction is resolved, the "next instruction" can execute.
- The hardware in each streaming multiprocessor dispatches an instruction each clock cycle if a ready instruction is available.
- The GPU in 11.4.5 supports 1024 such threads.
- What if our application needs more threads?
  - Threads are grouped into thread blocks.
  - Each thread block has up to 1024 threads (the HW limit).
  - The GPU can swap thread-block-in and out of main memory
    - This is GPU system software that we don't see as user-level programmers.

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 18 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 19 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 20 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 21 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 22 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 23 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 24 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 25 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 26 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 27 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 28 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 29 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 30 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 31 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 32 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 33 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 34 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 35 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 36 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 37 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 38 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 39 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 40 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 41 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 42 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 43 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 44 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 45 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 46 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 47 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 48 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 49 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 50 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 51 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 52 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 53 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 54 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 55 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 56 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 57 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 58 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 59 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 60 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 61 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 62 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 63 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 64 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 65 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 66 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 67 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 68 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 69 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 70 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 71 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 72 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 73 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 74 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 75 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 76 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 77 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 78 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 79 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 80 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 81 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 82 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 83 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 84 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 85 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 86 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 87 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 88 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 89 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 90 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 91 / 25

Greenstreet & Mitchell Introduction to GPGPUs CS 418 – Mar. 3, 2017 92 / 25

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 93 / 25

Greenstreet & Mitchell Introduction

- For this example, not really.
  - Execution time dominated by the memory copies.
- But, it shows the main pieces of a CUDA program.
- To get good performance:
  - We need to perform many operations for each value copied between memories.
  - We need to perform many operations in the GPU for each access to global memory.
  - We need enough threads to keep the GPU cores busy.
  - We need to watch out for **thread divergence**:
    - If different threads execute different paths on an if-then-else.
    - Then the else-threads stall while the then-threads execute, and vice-versa.
  - And many other constraints.
- GPUs are great if your problem matches the architecture.

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 23 / 25 Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 24 / 25 Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 25 / 25

## GPU Summary: architecture

- Lots of cores:
  - Up to 90 or more SIMD processors.
  - Each SIMD processor has 32 pipelines.
  - This is the NVIDIA architecture – other GPUs are similar.
- Deep, simple, execution pipelines
  - Optimized for floating point.
  - No bypassing: use multi-threading for performance.
  - Branches handled by predicated execution
    - "When you come to a fork in the road, take it."*  
(Often attributed to *Yogi Berra*)
- Limited on-chip memory
  - 1 or 2 MBs total. Big CPUs have 32-64MB of L3 cache.
  - The programmer manages data placement.

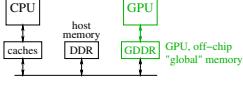
Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 2 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 3 / 24

## Data Parallelism

- When you see a **for-loop**:
  - Is the loop-index used as an array index?
  - Are the iterations independent?
  - If so, you probably have data-parallel code.
- Data-Parallel problems:
  - Run on GPU because each element (or segment) of the array can be handled by a different thread.
  - Data-parallel problems are good candidate for most parallel techniques because the available parallelism grows with the problem size.
  - Compare with "task parallelism" where the problem is divided into the same number of tasks regardless of its size.

Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 6 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 7 / 24

## Execution Model: Memory



- Host memory: DRAM and the CPU's caches
  - Accessible to host CPU but not to GPU.
- Device memory: GDDR DRAM on the graphics card.
  - Accessible by GPU.
  - The host can initiate transfers between host memory and device memory.
- The CUDA library includes functions to:
  - Allocate and free device memory.
  - Copy blocks between host and device memory.
  - BUT** host code can't read or write the device memory directly.

Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 10 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 11 / 24

## saxpy: host code (part 1 of 5)

```
int main(int argc, char **argv) {
 uint n = atoi(argv[1]);
 float *x, *y;
 float *dev_x, *dev_y;
 int size = n*sizeof(float);
 x = (float *)malloc(size);
 y = (float *)malloc(size);
 yy = (float *)malloc(size);
 for(int i = 0; i < n; i++) {
 x[i] = i;
 y[i] = i*i;
 }
}
```

- Declare variables for the arrays on the host and device.
- Allocate and initialize values in the host array.

Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 14 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 15 / 24

## saxpy: host code (part 5 of 5)

```
int main(void) {
 ...
 free(x);
 free(y);
 free(yy);
 cudaFree(dev_x);
 cudaFree(dev_y);
 exit(0);
}
```

- Clean up.
- We're done.

Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 18 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 19 / 24

## But it is fast?

- For the **saxpy** example as written here, not really.
  - Execution time dominated by the memory copies.
- But, it shows the main pieces of a CUDA program.
- To get good performance:
  - We need to perform many operations for each value copied between memories.
  - We need to perform many operations in the GPU for each access to global memory.
  - We need enough threads to keep the GPU cores busy.
  - We need to watch out for **thread divergence**:
    - If different threads execute different paths on an if-then-else.
    - Then the else-threads stall while the then-threads execute, and vice-versa.
  - And many other constraints.
- GPUs are great if your problem matches the architecture.

Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 22 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 23 / 24

March 6: Intro to CUDA  
Reading Kirk & Hwu 3rd ed., Ch. 2

March 8: CUDA Threads, Part 1  
Reading Kirk & Hwu Ch. 3

March 13: CUDA Threads, Part 2  
Reading Kirk & Hwu Ch. 4

March 15: CUDA Memory: examples  
Reading Kirk & Hwu Ch. 4

March 17: CUDA Performance  
Reading Kirk & Hwu Ch. 5

March 20: Matrix multiplication with CUDA, Part 1

March 22: Matrix multiplication with CUDA, Part 2

March 24 – April 3: Other Topics

more parallel algorithms, e.g. dynamic programming?

reasoning about concurrency, e.g. termination detection

other paradigms, e.g. Scala and futures?

April 5: Party: 50<sup>th</sup> Anniversary of Amdahl's Law

Mark Greenstreet

CpSc 418 – Mar. 6, 2017

- What is SIMD parallelism?
- How does a CUDA GPU handle branches?
- How does a CUDA GPU handle pipeline hazards?
- What is the difference between "shared memory" and "global memory" in CUDA programming?
- Think of a modification to the **saxpy** program and try it.
  - You'll probably find you're missing programming features for many things you'd like to try.
  - What do you need?
- Stay tuned for upcoming lectures.

GPU Summary: slide 2

CUDA

Data parallelism: slide 6

Program structure: slide 8

Memory: slide 10

A simple example: slide 12

Launching kernels: slide 19

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license <http://creativecommons.org/licenses/by/4.0/>

Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 25 / 25 Mark Greenstreet Introduction to GPGPUs CS 418 – Mar. 3, 2017 26 / 25

GPU Summary: Performance

- Today's processors are constrained by how much performance can be gotten using ~200 watts.
  - Moving bits around takes lots of energy.
  - $E \sim dt^{-\alpha}$ , where  $E$  is energy,  $d$  is distance,  $t$  is time per operation, and  $1 < \alpha < 2$  depending on design details.
  - Corollary:  $P \sim d^{\alpha-1}$ . Power grows somewhere between quadratically and cubically with clock frequency.
- How GPUs optimize performance/power
  - SIMD: instruction fetch and decode moves lots of bits. Amortize over many cores.
  - Simple pipelines: bypassing means moving bits quickly. GPUs omit bypasses.
  - High latency: avoid pipeline stages that must do a lot in a hurry.
  - Expose the memory hierarchy: let the programmer control moving data bits around.

Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 2 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 3 / 24

Which of the following loops are data parallel?

```
for(int i = 0; i < N; i++)
 c[i] = a[i] + b[i];

dotprod = 0.0;
for(int i = 0; i < N; i++)
 dotprod += a[i]*b[i];

for(int i = 1; i < N; i++)
 a[i] = 0.5*(a[i-1] + a[i]);

for(int i = 1; i < N; i++)
 a[i] = sqrt(a[i-1] + a[i]);

for(int i = 0; i < N; i++) {
 for(int j = 0; j < N; j++) {
 sum = 0.0;
 for(int k = 0; k < L; k++)
 sum += a[i,k]*b[k,j];
 c[i,j] = sum;
 }
}
```

Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 7 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 8 / 24

More Memory

- GPUs support fairly large off-chip memory bandwidth: 200-400GB/s.
  - But this isn't fast enough to keep 1000 processors busy at 1Gflop/s each!
- The GPU has on-chip memory to help:
  - Shared memory: 16KBs or 48KBs.
  - Registers: 128KBs (256KBs on more recent GPUs).
  - Note that we need to use each value from memory for 20 or more instructions or else the memory bandwidth will limit performance.
- GPUs also have L2 caches, around 1.5MBYTE in the most recent chips.
  - But I haven't found a good way to understand them from the textbook, or from other CUDA manuals.
  - The coherence/consistency guarantees seem to be pretty weak.

Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 12 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 13 / 24

Example: saxpy

saxpy = "Scalar a times x plus y".

The device code.

The host code.

The running saxpy

Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 16 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 17 / 24

saxpy: host code (part 3 of 5)

```
int main(void) {
 ...
 float a = 3.0;
 float x, y;
 float dev_x, dev_y;
 int size = n*sizeof(float);
 x = (float *)malloc(size);
 y = (float *)malloc(size);
 dev_x = (float *)malloc(size);
 dev_y = (float *)malloc(size);
 for(int i = 0; i < n; i++) {
 x[i] = i;
 y[i] = i*i;
 }
}
```

- Allocate arrays on the device.
- Copy data from host to device.

Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 18 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 19 / 24

Launching Kernels

- Terminology
  - Data parallel code that runs on the GPU is called a **kernel**.
  - Invoking a GPU kernel is called **launching** the kernel.
- How to launch a kernel
  - The host CPU invokes a **global** function.
  - The invocation needs to specify how many threads to create.
    - Example:
      - `a += c*ceil(n/256.0), 256>>(...)` says to create  $\lceil n/256 \rceil$  blocks of threads.
      - Each block consists of 256 threads.
      - See slide 20 for an explanation of threads and blocks.
      - The pointers to the arrays (in device memory) and the values of `n` and `c` are passed to the threads.

Copy the result back to the host.

Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 20 / 24 Mark Greenstreet Introduction to CUDA CS 418 – Mar. 6, 2017 21 / 24

Threads and Blocks

The GPU hardware combines threads into **warp**s

All of the threads of warp execute together – this is the SIMD part.

The functionality of a program doesn't depend on the warp details.

But understanding warps is critical for getting good performance.

Each warp has a "next instruction" pending execution.

If the dependencies for the next instruction are resolved, it can execute for all threads of the warp.

The hardware in each streaming multiprocessor dispatches an instruction each clock cycle if a ready instruction is available.

The GPU in Lin25 supports 32 such warps of 32 threads each in a "thread block".

If our application needs more threads?

Threads are grouped into "block blocks".

Each block has up to 1024 threads (the HW limit).

The GPU can swap thread-blocks in and out of main memory

\* This is GPU system software that we don't see as user-level programmers.

April 5: Party: 50<sup>th</sup> Anniversary of Amdahl's Law

Mark Greenstreet

CpSc 418 – Mar. 6, 2017

## GPU Summary: Economics

- GPUs are designed for the high-volume, consumer graphics market.
- Amitorize high design cost over a large number of units sold.
- This means GPUs aren't really optimized for scientific computing:
  - More on-chip memory would certainly help scientific computing, but not needed for graphics rendering.
  - Comparison: An nVidia GPU has about 2 MBytes of on-chip memory, an Intel Xeon can have 40MBytes or more.
  - Cache memory is about 60-70 transistors per byte.
  - A high-end nVidia GPU has 7 billion transistors, 1 or 2% for memory.
- Cheap is good
  - It's the economics of cheap computing that drives Moore's Law and all the other exponential growth-rate trends that make computing a field of intense, ongoing innovation.
  - That keeps the field in transition – deal with it.

## Programming GPUs: CUDA

Data Parallelism

CUDA program structure

Memory

Launching kernels

- A **global** function to called by the host program to execute on the GPU.
- There may be one or more **device** functions as well.

- One or more host functions, including **main** to run on the host CPU.
  - Allocate device memory.
  - Copy data from host memory to device memory.
  - "Launch" the device kernel by calling the **global** function.
  - Copy the result from device memory to host memory.

A **global** function to called by the host program to execute on the GPU.

The device code.

The host code.

The running saxpy

## saxpy: device code

```
__global__ void saxpy(uint n, float a, float *x, float *y) {
 uint i = blockIdx.x*blockDim.x + threadIdx.x;
 if(i < n)
 y[i] = ax[i] + y[i];
}
```

- Each thread has `x` and `y` indices.
- We'll just use `x` for this simple example.

- Note that we are creating one thread per vector element:
- Exploits GPU hardware support for multithreading.
- We need to keep in mind that there are a large, but limited number of threads available.

## saxpy: host code (part 4 of 5)

```
...
for(int i = 0; i < n; i++) // check the result
 if(y[i] != a*x[i] + y[i]) {
 fprintf(stderr,
 "ERROR: i=%d, a[%d]=%f, b[%d]=%f, c[%d]=%f\n",
 i, a[i], b[i], c[i]);
 exit(-1);
 }
printf("The results match!\n");
}

Check the results.
```

## Compiling and running

lin25\$ nvcc saxpy.cu -o saxpy

lin25\$ ./saxpy 1000

The results match!

## CUDA Threads

Mark Greenstreet & Ian M. Mitchell

CPSC 418 – March 8 & March 10, 2017

- Kernel organization: grids, blocks & threads.

- Hardware organization: SMs, SPs & warps.

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet & Ian M. Mitchell and are made available under the terms of the Creative Commons Attribution 4.0 International license <http://creativecommons.org/licenses/by/4.0/> and are made available under the terms of the Creative Commons Attribution 4.0 International license <http://creativecommons.org/licenses/by/4.0/>

## Compute Capability

- Lots of Nvidia jargon here.
- Lots of very specific constraints on hardware capabilities.
- Values of those constraints depend on the compute capability: essentially a version number for the GPU hardware.
  - CS department lab (lin01, lin02, ..., lin23) has a `gridres.c` which has GeForce GTX 550 Ti which features compute capability 2.1.
  - Examples of recent GPUs:
    - Compute capability 3.5: GTX 730 & GTX 780.
    - Compute capability 5.0: GTX 750, 800 & 960M.
    - Compute capability 5.2: GTX 9x, 950M.
    - Compute capability 6.1: GTX 10xx.
  - More details at the CUDA wikipedia page.

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 2 / 20

### Threads and blocks: launching a kernel

- Let's say we have:
 

```
global void kernelFun(args)
```
- To launch this kernel, we execute a statement like:
 

```
kernelFun<<<dimGrid, dimBlock>>>(actuals);
```
- where
  - `dimGrid` specifies the dimension(s) of the grid (an array of blocks):
    - `dimGrid` can be an `int`, in which case the array is 1D.
    - `dimGrid` can be a `dim3`, for example:
 

```
dim3(6, 4, 1);
```
  - `dimBlock` specifies the dimension(s) of each block (an array of threads):
    - `dimBlock` can be an `int` or a `dim3`.

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 6 / 20

### Bounds checking: in the kernel

- The kernel launch looks like:
 

```
kernelFun<<<ceil(n/256.0), 256>>>(n, myArray);
```
- THINK:** what if `n` is not a multiple of 256?
  - We'll launch more than `n` threads?
  - For example, if `n==1000`, then we'll launch 4 blocks of 256 threads. A total of 1024 threads.
  - What will the last 24 threads do?
- Add a test:
 

```
uint myIdx = blockDim.x*blockIdx.x + threadIdx.x;
if(myIdx < n) {
 ...
}
```

## A Warped Example: Reduce (part 1)

- Consider a reduce of an array `data` containing `n` elements using `n/2` threads (assume `n` is power of 2).
- Simple code:
 

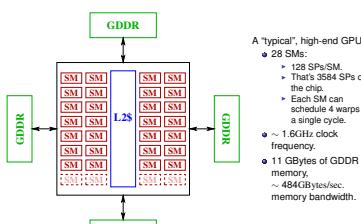
```
for(int stride = 1; stride < n; stride += stride) {
 if(myIdx & (stride-1) == 0)
 data[2*myIdx] += data[2*myIdx + stride];
 __syncthreads();
}
```
- The `__syncthreads()` call ensures that every thread has completed an iteration of the loop before any thread starts the next iteration.
  - More discussion on [slide 18](#).

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 14 / 20

### Synchronization

- The reduce example used `__syncthreads()`: all the threads in the block must execute this statement before any can continue beyond it.
  - Be very careful about thread divergence: All threads in the block must meet at the same barrier.
  - That means the **same** line of code.
  - In loops, that means the **same** iteration.
  - Executing different `__syncthreads()` commands will cause the kernel to hang.
- Also, `__syncthreads()` only synchronizes between threads within a single block.
  - Note that threads within a warp already stay synchronized because they are executed together.
  - The only way to synchronize between threads in different blocks is to finish the kernel and launch another.
- We'll cover synchronization in more detail later.

## First, GPU Architecture Review



## Greenstreet & Mitchell CUDA Memory CS/SC 418 - Mar. 13 & 15, 2017 2 / 24

- Each SP has its own register file.
- The register file is partitioned between threads executing on the SP.
- Local variables are placed in registers.
  - The compiler in-lins functions when it can
    - A kernel with recursive functions or deeply nested calls can cause register spills to main memory – this is **slow**.
  - Local array variables are mapped to global memory – **watch out**.

## Thread organization: Grids, Blocks and Threads

- When a kernel is launched, it creates a collection of threads.
- This collection is called a **grid**.
  - A grid is organized as an array of **blocks**.
  - Each block is an array of **threads**.
  - Array sizes are fixed once a kernel is launched.
- Why so many details?
  - Switching between blocks is done (infer) by software in the GPU.
  - Switching between threads in a block is done by hardware.
  - By distinguishing blocks from threads, the CUDA model exposes the performance issues to the programmer.

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 3 / 20

### Threads and Blocks within a Kernel's Grid

- Within a running kernel, CUDA-C provides four built-in variables to determine the position of a thread within the grid: `blockDim`, `blockIdx`, `threadDim`, and `threadIdx`.
- There is a naming pattern:
  - Each of these structures has three fields: `x`, `y` and `z` corresponding to the three possible dimensions.
  - `blockIdx.x` gives the size of the grid in each dimension `x`, `y` or `z`.
  - `threadIdx.x` gives the size of each block in each dimension.
  - `blockIdx.z` gives the indices of the thread's block within the grid.
  - `threadIdx.z` gives the indices of the thread within its block.
- For dimensions which are absent:
  - `blockDim.x` or `threadDim.x` will be 1.
  - `blockIdx.x` or `threadIdx.x` will be 0.

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 6 / 20

### SMs, SPs and Warps (oh my!)

- Each **streaming multiprocessor** (SM) has multiple **streaming processors** (SPs) and can be responsible for multiple groups of 32 threads called **warps**.
  - From the *New Oxford American Dictionary*: (the) "warp" is "the threads on a loom over and under which other threads (the weft) are passed to make cloth"
- Details, details...
  - These concepts are not part of the CUDA platform and API: Code is written in terms of a grid of blocks of threads.
  - You can write correct code without thinking about these details.
  - If you want to write fast code, you must take them into account.
  - The block vs grid structure exposes these details if you want to take advantage of them.

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 7 / 20

### SMs, SPs and Warps: What are They?

- Each streaming multiprocessor (SM) in the GPU executes threads in SIMD fashion.
  - All threads in a block are assigned to the same SM.
  - Each SM has a single (or small number of?) instruction fetch unit(s) and a large number of execution units.
- Each SM has multiple streaming processors (SPs) that actually execute an instruction.
  - The SPs are specialized: ALU, load / store, special function units.
  - A single SP can perform a single operation on a small set of threads.
- A warp is a collection of 32 threads that execute together on the same SP.

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 8 / 20

### SMs, SPs and Warps: Where are We?

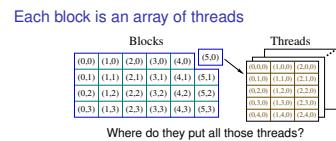
- Note the constraints:
 

```
0 ≤ blockIdx.x < blockDim.x
0 ≤ blockIdx.y < blockDim.y
0 ≤ blockIdx.z < blockDim.z
0 ≤ threadIdx.x < threadDim.x
0 ≤ threadIdx.y < threadDim.y
0 ≤ threadIdx.z < threadDim.z
```
- Because the size of blocks are limited, it is common to use code such as:
 

```
uint myIdx = blockDim.x*blockIdx.x + threadIdx.x;
```

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 9 / 20

### Bounds checking: launching kernels



Where do they all put those threads?

- Threads are scheduled by the GPU **hardware**.
- Blocks can be arranged as 1D, 2D or 3D array
  - Dimensions are called "`x`", "`y`" and "`z`".
- There can be lots of blocks:
  - Each dimension can be up to  $2^{16} - 1 = 65535$ .
  - CC 3.0+ allows `x` dimension up to  $2^{31} - 1$  blocks.
- However, total number of threads per block (product of all dimensions) is also capped at 1024.

## A Warped Example: Reduce (part 2)

- Consider a reduce of an array `data` containing `n` elements using `n/2` threads (assume `n` is power of 2).
- Simple code:
 

```
for(int stride = 1; stride < n; stride += stride) {
 if(myIdx & (stride-1) == 0)
 data[2*myIdx] += data[2*myIdx + stride];
 __syncthreads();
}
```
- The `__syncthreads()` call ensures that every thread has completed an iteration of the loop before any thread starts the next iteration.
  - More discussion on [slide 18](#).

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 10 / 20

### A Warped Example: Reduce (part 2)

- Consider `n == 16`
  - First iteration, for `i` in 0, ..., 7:
 

```
data[2*i] += data[2*i+1]
```
  - Now, all the even indexed elements have their sum with their odd counterpart.
  - Second iteration, for `i` in 0, 2, 4, 6:
 

```
data[2*i] += data[2*i+2]
```
  - All elements with indices that are multiples of four, have their sum with the next three elements.
  - Third iteration leads with `data[0]` and `data[8]` holding sums for their halves of the array.
  - The fourth iteration puts the complete sum into `data[0]`.
  - There are at most 8 threads working, so everything fits within a single warp.

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 11 / 20

### A Warped Example: Reduce (part 3)

- What if `n==1024`?
  - We have 512 threads: 16 warps of 32 threads.
  - In the first iteration all threads are active.
  - In the second iteration each warp has 16 active threads, so the GPU has to execute the code for all 16 warps even though half the threads do nothing.
  - In subsequent iterations, the warps are more and more poorly utilized.
- This solution is correct, but much of the parallel hardware will sit idle much of the time.
- We would like to pack the busy threads into the minimum number of warps.

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 12 / 20

### Warp Speed!

```
for(int stride = n/2; stride > 0; stride >>= 1) {
 if(myIdx < stride)
 data[myIdx] += data[myIdx + stride];
 __syncthreads();
}
```

- Consider `n == 1024` again.
  - In the first iteration, there are 16 active warps – all threads in each warp are busy.
  - In the second iteration, there are 8 active warps – all threads in each active warp are busy.
  - Similarly, for the 3<sup>rd</sup> through 5<sup>th</sup> iterations
- The number of active warps decreases after each iteration, but all threads in each active warp are busy.
  - The inactive warps have no pending instructions, so they will not be scheduled and will not occupy processing resources.

## Greenstreet & Mitchell CUDA Threads CS/SC 418 - Mar. 8 & 10, 2017 13 / 20

### Review

## Greenstreet & Mitchell CUDA Memory CS/SC 418 - March 13 & 15, 2017 1 / 24

- In CUDA, what is a grid, a block, and a thread?
- Why does CUDA allow millions of thread blocks but only 1024 threads per block?
- How does a programmer specify the number of blocks and number of threads when launching a CUDA kernel?
- How does a thread determine its position within the grid?
- Why do threads need to check their indices against array bounds?
- What is a warp? Why does it matter?

## CUDA: Memory

Mark Greenstreet

CpSc 418 – March 13 & 15, 2017

## Greenstreet & Mitchell CUDA Memory CS/SC 418 - Mar. 13 & 15, 2017 17 / 24

### Why do we need a memory hierarchy

- Architecture Snapshot
  - Registers
  - Shared Memory
  - Global Memory
  - Other Memory: texture memory, constant memory, caches
  - Summary, preview, review, tide-chart
- Unless otherwise noted or obvi, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license
 <http://creativecommons.org/licenses/by-nd/4.0/>
- Focus on the innermost loop: `for (k ...)`
  - Why?
  - How many floating point operations per iteration?
  - How many memory reads?
  - How many memory writes?
  - What is the "Compute-to-Global-Memory-Access" ratio (CGMA)?

## Greenstreet & Mitchell CUDA Memory CS/SC 418 - Mar. 13 & 15, 2017 18 / 20

### A Streaming Multiprocessor (SM)

- A "typical", high-end GPU.
  - 28 SMs:
    - Each SM has 32 SPs on the chip.
    - Each SM can support up to 4 warps in a single cycle.
    - ~1.6GHz clock frequency.
    - ~11 GBytes of GDDR memory, ~484GB/sec., memory bandwidth.
- Each of the pipelines is an SP (streaming processor)
  - Lots of deep pipelines.
  - Lots of threads: when we encounter an architectural challenge:
    - Raising throughput is **easy**, lowering latency is **hard**.
    - Solve problem by increasing latency and adding threads.
  - Make the programmer deal with it.

## Greenstreet & Mitchell CUDA Memory CS/SC 418 - Mar. 13 & 15, 2017 19 / 20

### Registers and Memory Bandwidth

- The GPU on [slide 2](#) has 28 SMs, each with 128 SPs.
  - Each SP has access to a register file.
  - I'll guess two register reads and one write per clock cycle, per SP.
  - I'll assume 4-byte registers.
  - We get
 
$$28SM \times 128SP \times 3RW \times \frac{1.6 \text{ GHz}}{\text{SP} \times \text{cycle}} \times \frac{4 \text{ Byte}}{\text{sec.}} = 68613 \frac{\text{GByte}}{\text{sec.}}$$
  - 142 times faster than main memory bandwidth!
  - ymmv: the GPUs in the linXX box are older.
- Registers
  - Each SM has 256K registers, and 64 active warps, with 32 threads/warp.
    - That's 32 4-byte registers per thread.
    - If a thread uses more registers
      - The SM cannot fully use its warp scheduler, or Registers will spill to main memory — **slow**
    - The numbers are smaller for older GPUs.
      - The GTX 550 Ti GPUs in the linXX boxes support 21 registers/thread.

## Greenstreet & Mitchell CUDA Memory CS/SC 418 - Mar. 13 & 15, 2017 20 / 20

### Registers and Thread Scheduling

- Registers
  - Each SP has its own register file.
  - The register file is partitioned between threads executing on the SP.
  - Local variables are placed in registers.
    - The compiler in-lins functions when it can
      - A kernel with recursive functions or deeply nested calls can cause register spills to main memory – this is **slow**.
  - Register spills are mapped to global memory – **watch out**.
- Registers and Memory Bandwidth
  - The GPU on [slide 2](#) has 28 SMs, each with 128 SPs.
    - Each SM has access to a register file.
    - I'll guess two register reads and one write per clock cycle, per SP.
    - I'll assume 4-byte registers.
    - We get
 
$$28SM \times 128SP \times 3RW \times \frac{1.6 \text{ GHz}}{\text{SP} \times \text{cycle}} \times \frac{4 \text{ Byte}}{\text{sec.}} = 68613 \frac{\text{GByte}}{\text{sec.}}$$
    - 142 times faster than main memory bandwidth!
    - ymmv: the GPUs in the linXX box are older.

## Greenstreet & Mitchell CUDA Memory CS/SC 418 - Mar. 13 & 15, 2017 21 / 24

### Registers and Thread Scheduling

- Registers and Memory Bandwidth
  - The GPU on [slide 2](#) has 28 SMs, each with 128 SPs.
    - Each SM has access to a register file.
    - I'll guess two register reads and one write per clock cycle, per SP.
    - I'll assume 4-byte registers.
    - We get
 
$$28SM \times 128SP \times 3RW \times \frac{1.6 \text{ GHz}}{\text{SP} \times \text{cycle}} \times \frac{4 \text{ Byte}}{\text{sec.}} = 68613 \frac{\text{GByte}}{\text{sec.}}$$
    - 142 times faster than main memory bandwidth!
    - ymmv: the GPUs in the linXX box are older.

## Greenstreet & Mitchell CUDA Memory CS/SC 418 - Mar. 13 & 15, 2017 22 / 24

### Registers and Thread Scheduling

- Registers and Memory Bandwidth
  - The GPU on [slide 2](#) has 28 SMs, each with 128 SPs.
    - Each SM has access to a register file.
    - I'll guess two register reads and one write per clock cycle, per SP.
    - I'll assume 4-byte registers.
    - We get
 
$$28SM \times 128SP \times 3RW \times \frac{1.6 \text{ GHz}}{\text{SP} \times \text{cycle}} \times \frac{4 \text{ Byte}}{\text{sec.}} = 68613 \frac{\text{GByte}}{\text{sec.}}$$
    - 142 times faster than main memory bandwidth!
    - ymmv: the GPUs in the linXX box are older.

## Greenstreet & Mitchell CUDA Memory CS/SC 418 - Mar. 13 & 15, 2017 23 / 24

### Registers and Thread Scheduling

- Registers and Memory Bandwidth
  - The GPU on [slide 2](#) has 28 SMs, each with 128 SPs.
    - Each SM has access to a register file.
    - I'll guess two register reads and one write per clock cycle, per SP.
    - I'll assume 4-byte registers.
    - We get
 
$$28SM \times 128SP \times 3RW \times \frac{1.6 \text{ GHz}}{\text{SP} \times \text{cycle}} \times \frac{4 \text{ Byte}}{\text{sec.}} = 68613 \frac{\text{GByte}}{\text{sec.}}$$
    - 142 times faster than main memory bandwidth!
    - ymmv: the GPUs in the linXX box are older.

## Greenstreet & Mitchell CUDA Memory CS/SC 418 - Mar. 13 & 15, 2017 24 / 24

### Registers and Thread Scheduling

- Registers and Memory Bandwidth
  - The GPU on [slide 2](#) has 28 SMs, each with 128 SPs.
    - Each SM has access to a register file.
    - I'll guess two register reads and one write per clock cycle, per SP.
    - I'll assume 4-byte registers.
    - We get
 
$$28SM \times 128SP \times 3RW \times \frac{1.6 \text{ GHz}}{\text{SP} \times \text{cycle}} \times \frac{4 \text{ Byte}}{\text{sec.}} = 68613 \frac{\text{GByte}}{\text{sec.}}$$
    - 142 times faster than main memory bandwidth!
    - ymmv: the GPUs in the linXX box are older.

## Registers and Matrix Multiply

```
for(int i = 0; i < M; i += 2) {
 for(int j = 0; j < N; j += 2) {
 sum0 = sum1 = sum10 = sum11 = 0.0;
 for(int k = 0; k < L; k += 2) {
 a00 = a[i,k]; a01 = a[i,k+1];
 a10 = a[i+1,k]; a11 = a[i+1,k+1];
 b00 = b[k,j]; b01 = b[k+1,j];
 b10 = b[k,j]; b11 = b[k+1,j+1];
 sum0 += a00*b00 + a01*b01;
 sum1 += a00*b10 + a01*b11;
 sum10 += a10*b00 + a11*b01;
 sum11 += a10*b10 + a11*b11;
 }
 c11[j] = sum0; c11[j+1] = sum1;
 c11[j+2] = sum10; c11[j+3] = sum11;
 }
}
```

- use a register to accumulate  $c[i,j]$
- hold blocks of each matrix in main memory.
- can use each value loaded from  $a[i,k]$  or  $b[k,j]$  three or four times.
- What is the CGMA for the example above?

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 10 / 24

## Shared Memory: Collisions



- When one thread in a warp accesses shared memory, all active threads in the warp access shared memory.
- If each thread accesses a different bank, then all accesses are performed in a single cycle.
  - Otherwise, the load or store can take multiple cycles.
  - Multiple accesses to the same bank are called **collisions**.
  - The **worst-case** occurs when **all threads access the same bank**.
- The programmer needs to think about the index calculations to avoid collisions.
- When programming GPUs, the programmer needs to think about index calculations a lot.

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 14 / 24

## Writing and reading DRAM

- Writing:** easy
  - drive all 1024 column-lines to the values you want to write.
  - open up all the rows for one row.
  - the row is holding caps for each column in that row get filled or emptied.
  - note: you end up writing **every column** in the row; so writes are often preceded by reads.
- Reading:** hard
  - drive all 1024 column-lines to "half-way", and let them float\*.
  - open up all the values for one row.
  - if the level in the pipe goes up a tiny amount, that cup held a 1.
  - if the level in the pipe goes down a tiny amount, that cup held a 0.
  - it's a delicate measurement – it takes time to set it up.
  - This is why DRAM is slow.
- But: we just read 1024 bits, from each chip of the DIMM.
  - That's 16GBs = 2bytes total.
  - Conclusion: DRAM has awful latency, but we can get very good bandwidth.
    - bandwidth bottleneck is the wires from the DIMM to the CPU or GPU.
    - But I'm pretty sure that Ian won't let me give a lecture on transmission lines, phase-locked loops, equalizers, etc. and the other cool stuff in the DDR (or RDIMM) interface.

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 18 / 24

## Summary

- GPUs can have thousands of execution units, but only a few off-chip memory interfaces.
  - This means that the GPU can perform 10-50 floating point operations for every memory read or write.
  - Arithmetic operations are very cheap compared with memory operations.
- To mitigate the off-chip memory bottleneck
  - GPUs have, limited on-chip memory
  - Registers and the per-block, shared-memory will be our main concern in this class.
- Moving data between different kinds of storage is the programmer's responsibility.
  - The programmer explicitly declares variables to be stored in shared memory.
  - The programmer needs to be aware of the per-thread register usage to achieve good SM utilization.
  - The only way to communicate between thread blocks is to write to global memory, end the kernel, and start a new kernel (ouch!)

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 22 / 24

## CUDA: Performance Considerations

Mark Greenstreet & Ian M. Mitchell

CPSC 418 – March 17, 2017

- [Thread Divergence](#)
- [Floating Point Follies](#)
- [Memory Accesses](#)
- [Occupancy](#)
- [Granularity](#)

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet & Ian M. Mitchell  
http://csweb.cs.uvic.ca/~mitchell/courses/cs418/4.0/International License  
Attribution 4.0 International license

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 1 / 22

## Fused multiply adds

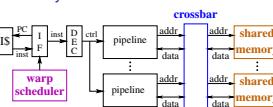
- Calculating  $ax + b$  is very common
  - Example: dot product.
- The multiplier hardware is just a pipeline of adders.
  - When multiplying  $a \times x$ , the hardware can start the pipeline from  $b$  instead of from  $0$ .
  - We get the sum for "free".
  - This is called a **fused** multiply-add.
- The marketing people like to count the fused multiply-add as **two** floating point operations.
  - This helps make some performance claims make sense.
- For the obsessive compulsive:
  - Rounding with a fused multiply add can be slightly different than when doing two separate operations.
  - Compilers usually let the users specify "strict" floating point (no fusing) or "fast" floating point (with fusing).
  - `NVCC` uses fused multiply add unless you give it an option not to.

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 5 / 22

- ### Memory System Parallelism 3: Independent Memory Components
- Even after memory address is delivered, it still takes time for the DRAM to return the data.
- Rather than let the memory bus sit idle while waiting, pipeline a bunch of memory requests to different memory components.
    - K&H(3) calls these "banks".
    - Mark called these "tiles".
  - Consecutive memory chunks are assigned first to channels / banks (see previous slide).
    - These subsystems allow concurrent access because they have independent communication lines.
  - Then assign next set of consecutive chunks to banks / tiles.
    - These subsystems allow sequential but pipelined access because they share communication lines.
- Pipelining increases throughput (although latency remains).
- Only relevant for global memory.
  - Shared memory achieves dramatically lower latency with SRAM.

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 9 / 22

## Shared Memory



- On-chip, one bank per SP
- Banks are interleaved by:
  - Early CUDA GPUs: 4-byte word
  - Later GPUs: programmer configurable 4-byte or 8-byte words
  - Why?
- Shared memory is a limited resource: 48KBytes to 96Kbytes/SM.
  - Each SM has more registers than shared-memory.
  - Shared memory demands limit how many blocks can execute concurrently on a SM.

See notes on the next slide.

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 11 / 24

## Shared Memory Example: Reduce

## Shared Memory Example: Reduce

```
shared_ float v[1024];
device_ float t(float x) {
 return((5/2)*x - x);
}
device_ void compute_and_reduce(uint n, uint m, float *x) {
 myId = threadIdx.x;
 if(myId < n) {
 float y = x[myId];
 for(uint i = 0; i < m; i++) {
 y = t(y);
 y = v[i] + y;
 }
 v[myId] = y;
 }
 __syncthreads();
 if(myId < m) {
 v[myId] += v[myId+n];
 }
 x[myId] = v[myId];
}
```

See notes on the next slide.

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 11 / 24

## Shared Memory Example: Matrix Multiply

- Running example from the textbook:  $C = AB$
- Each thread-block loads a  $16 \times 16$  block from  $A$  and  $B$ .
    - The threads to these loads "cooperatively":
      - Read  $A_{i,j}$  and  $B_{i,j}$  from global memory with "coalesced" loads.
      - Write these blocks to shared-memory in a way that avoids bank conflicts.
  - Compute:  $C_{i,j} = A_{i,j} \times B_{i,j}$ .
    - This takes  $16^2 = 4096$  fused multiply-adds.
    - Loading  $A_{i,j}$  fetches  $16^2 = 256$  floats from global memory.
    - Likewise for  $B_{i,j}$ . Total of 512 blocks fetched.
    - CGMA =  $4096/512 = 8$ .
  - Note: The L2 cache might help here:  $A$  and  $B$  are read-only.
    - Need to try more experiments.

## DRAM meet DRAM

- DRAM summary: terrible latency (60-200ns or more), fairly bandwidth.
  - The GPU lets the program take advantage of high bandwidth:
    - If the 32 bards from a warp access 32 consecutive memory location,
      - The GPU does **one** GDDR access,
      - and it transfers a large block of data.
    - The same optimization is applied to stores, and to loads from the on-chip caches.
- In CUDA-speak, if the loads from a warp access consecutive locations, we say that the memory accesses are **coalesced**.
  - Note that the memory optimizations are exposed to the programmer.
  - You can get the performance by considering the memory model.
  - But, it's not automatic.

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 15 / 24

## Global Memory

## Example: Matrix Multiplication

- In C, matrices are usually stored in row-major order.
  - $A[i,k]$  and  $A[i,k+1]$  are at adjacent locations, but  $B[k,j]$  and  $B[k+1,j]$  are  $N$  words apart (for  $N \times N$  matrices).
- For matrix multiplication, accesses to  $A$  are naturally coalesced, but accesses to  $B$ .
  - The optimized code loads a block of  $B$  into shared memory.
    - This allows access to be coalesced.
    - But we need to be careful about how we store the data in the shared memory to avoid bank conflicts.

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 16 / 24

## Other Memory

## Notes on Reduce

- We calculate  $\sum(x[\text{myId}])$  locally using the register variable  $y$ .
- For the reduce, we use the **shared array**  $v$ .
  - This avoids the penalty of off-chip, global memory access for each step of the reduce.
  - All threads in a warp can access shared memory on the same cycle.
  - We can have multiple blocks running on multiple SMs
    - Each SM has its own shared memory.
    - Blocks running on different SMs in parallel can **all** access their shared memories in parallel.
    - But**, threads in one block do not share shared-memory with threads in other blocks.
  - To perform a reduce across blocks:
    - Each block writes its subtotal to the **global** memory.
    - The results from the blocks are combined on the host CPU or by launching a new kernel.
- At the end, we copy our value from shared memory to the global memory so the CPU or a subsequent kernel can access it.

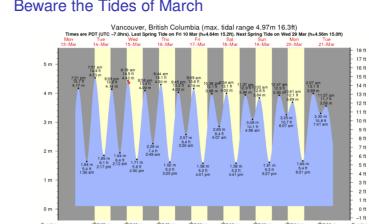
## Now for a word about DRAM

The memory that you plug into your computer is mounted on DIMMs (dual-inline memory modules).

- A DIMM typically has 16 or 18 chips
- E.g. each chip of an 8Gbyte DIMM holds 512MBytes = 4Gbits.
- Each chip consists of many "tiles",
  - a typical chip has 1Mbit/tile
  - that's 4096 tiles for a 4Gb chip.
- Each tile is an array of capacitors.
  - each capacitor holds 1 bit.
  - a typical tile could have 1024 rows and 1024 columns.

- Constant memory: cached, read-only access of global memory.
- Texture memory: global memory with special access operations.
- L1 and L2 caches: only for memory reads.

## Beware the Tides of March



From <https://www.tide-forecast.com/locations/Vancouver/British-Columbia/tides/tides.html>

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 21 / 24

## Review

- What is CGMA?
- On slide 15 we computed the CGMA for matrix-multiplication using 16  $\times$  16 blocks of the  $A$ ,  $B$ , and  $C$  matrices.
  - How many such thread-blocks can execute concurrently on an SM with 48KBytes of memory?
  - How does the CGMA change if we use 32  $\times$  32 blocks?
  - If we use the larger matrix-blocks, how many thread-blocks can execute concurrently on an SM with 48KBytes of memory?
  - If we use the larger matrix-blocks, how many thread-blocks can execute concurrently on an SM with 96Kbytes of memory?
- What are bank conflicts?
- How can increasing the number of registers used by a thread improve performance?
- How can increasing the number of registers used by a thread degrade performance?
- What is a "coalesced memory access"?

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 23 / 24

## Review

- When working on my solution to [last year's HW3](#), OI,
  - I first wrote:
 

```
x = alpha*x*(1.0 - x);
```
  - and the performance was disappointing.
  - After many frustrating attempts to track down the problem, I added one, little f:
 

```
x = alpha*x*(1.0f - x);
```
  - and my code ran 5.5x faster.

What happened?

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 24 / 24

## Remarks about floating point

- When working on my solution to [last year's HW3](#), OI,
  - I first wrote:
 

```
x = alpha*x*(1.0 - x);
```
  - and the performance was disappointing.
  - After many frustrating attempts to track down the problem, I added one, little f:
 

```
x = alpha*x*(1.0f - x);
```
  - and my code ran 5.5x faster.

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 25 / 24

## Floating point

- When working on my solution to [last year's HW3](#), OI,
  - I first wrote:
 

```
x = alpha*x*(1.0 - x);
```
  - and the performance was disappointing.
  - After many frustrating attempts to track down the problem, I added one, little f:
 

```
x = alpha*x*(1.0f - x);
```
  - and my code ran 5.5x faster.

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 25 / 24

## Memory System Parallelism 1: Get Lots of Data

Memory is addressed per byte, but you retrieve a bunch of (sequential) bytes at once.

- GDDR5 DRAM: 32-bit bus per chip and transfers are in 16 word bursts (so 64 bytes per access per chip).
- GPU global memory (from GDDR DRAMs): Accessed by 32-, 64- or 128-byte transactions.
  - Transactions must be "naturally" aligned: First address must be a multiple of the transaction size.
  - CC 2.x: L1 cache (1 per SM) serviced by 128-byte transactions, L2 cache (shared by SMs) by 32-byte transactions.
  - CC 6.x: Same as 2.x, but L1 cache rules are complicated.
- GPU shared memory (on-chip SRAM): Access in 32-bit words.

Amortize addressing overhead and thereby increase bandwidth.

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 3 / 22

## Memory System Parallelism 2: Multiple Interfaces

If one memory component cannot give you enough bandwidth, use a backup (see [March 13 slides](#)).

- Global memory: K&H(3) calls these "channels" (March 13 slide 2).
- Shared memory: Mark called these "banks" (March 13 slide 11) and NVIDIA documentation does too.
- Do not confuse with K&H(3)'s "banks" (see next slide).

Consecutive chunks are placed into components in a round-robin fashion, where "chunk" means

- 32-bytes (64 more recently?) in global memory.
- 32- or 64-bits in shared memory.

Separate subsystems can all provide data at their native rate and thereby increase bandwidth.

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 4 / 22

## Implications for Shared Memory

Try to get memory access addresses from threads in a warp to be very close together.

- Accesses to consecutive (or nearly so) addressees are coalesced into a single transaction on the off-chip memory bus.

You should already be doing this for your CPU designs so that your caches can take advantage of spatial locality.

- Best coalescing occurs when the set of addresses is naturally aligned.

For two and higher dimensional arrays, that may mean padding thread block and array width allocation in memory to be a multiple of the warp size.

- Possibility of channel / bank collisions would argue for avoiding addresses with the same "middle" bits.

I could not find NVidia documentation of these details.

How do caches interact with channels / banks?

Comments from Mark?

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 12 / 22

## Implications for Global Memory

Try to get memory access addresses from threads in a warp to be very close together.

- Accesses to consecutive (or nearly so) addressees are coalesced into a single transaction on the off-chip memory bus.

You should already be doing this for your CPU designs so that your caches can take advantage of spatial locality.

- Best coalescing occurs when the set of addresses is naturally aligned.

For two and higher dimensional arrays, that may mean padding thread block and array width allocation in memory to be a multiple of the warp size.

- Possibility of channel / bank collisions would argue for avoiding addresses with the same "middle" bits.

I could not find NVidia documentation of these details.

How do caches interact with channels / banks?

Comments from Mark?

Mark Greenstreet CUDA Memory CS 418 – Mar. 13 & 15, 2017 12 / 22

## SMs and Thread Occupancy

- Occupancy: how many warps are available for the SM
  - Why we care: the SP pipelines have long latencies.
  - The CUDA approach is to run lots of threads simultaneously to keep the pipelines busy.
- Limits to occupancy
  - How many blocks per SM.
  - How much shared-memory per block.
  - How many threads per block.
  - How many registers per thread.
- Figuring it out
  - nvcc -O3 -c -ptxas-options -v examples.cu
  - The nVidia occupancy calculator: [CUDA\\_Occupancy\\_calculator.xls](#)
  - But we can do it manually?

## Occupancy with CUDA 2.1

- Different GPUs at level CUDA 2.1 have differing numbers of SMs.
  - But the SMs all look the same, even for different GPUs.
- CUDA 2.1 SMs
  - An SM has warps of 32 threads
  - An SM can simultaneously execute up to 1536 threads (48 warps).
  - An SM has 32K (2<sup>15</sup>) 32-bit registers (128Kbytes, 1K registers/SP).
  - An SM has 48K bytes of shared memory.
  - An SM can simultaneously execute up to 8 blocks.
  - Each block can have up to 1024 threads.

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 13 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 14 / 22

## Why all these numbers?

- When designing a new generation of GPUs, the GPU architects run lots of simulations to estimate the performance for various choices of the architectural parameters.
- For example, if more warps are allowed in the scheduling pool
  - The SM will have useful instructions to dispatch more often → better performance.
  - **BUT** the on-chip circuitry to hold and manage the scheduling pool will be larger.
  - This means instruction scheduling will be slower → a longer clock period.
  - Instruction scheduling will use more power → a longer clock period, or fewer SMs, or more expensive chip cooling.
  - The real-state on the chip could have been used for something else. Is this the best use of that area?
  - Note that CUDA 5 made the increase to 64 warps/SM.
- Architects explore these trade-offs to optimize performance for graphics applications, the main source of revenue.
- Architects are also risk-adverse: make the chip as much like the last one that worked as you can.
- These hard-wired constraints have a large impact on program performance.

## SMs, blocks, and threads

- A SM can simultaneously execute most 8 blocks.
- All blocks have the same number of threads.
- Thus, a SM can execute at most

$$\min \left( 8, \left\lfloor \frac{1536}{\text{threadsPerBlock}} \right\rfloor \right)$$

blocks.

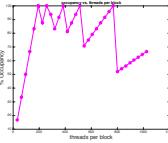
- The ratio of the number of threads executing to the maximum possible is called the "thread occupancy".

$$\text{threadOccupancy} \leq \frac{\min \left( 8, \left\lfloor \frac{1536}{\text{threadsPerBlock}} \right\rfloor \right) \text{threadsPerBlock}}{1536}$$

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 13 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 14 / 22

## SMs, blocks, and threads – the plot



- I get 100% occupancy when  $\text{threadsPerBlock} \in \{192, 384, 768\}$ , but the CUDA calculator doesn't.
  - I'll have to try some experiments – stay tuned.
- This assumes the grid had enough blocks to keep the SMs busy.
  - A grid with a single block will have poor performance.

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 17 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 18 / 22

## Granularity

### How much work should a kernel do?

- Do more work within a kernel: Launching each kernel takes time.
- Do less work within a kernel: New kernels allow for changes in block and grid size, and ensure synchronization between threads even in different blocks.
- Either way: Minimize movement of data to and from the host.
- How much work should a thread do?
  - Do more work in a single thread: Fewer chances for memory collisions, easier synchronization, less register contention.
  - Do less work in a single thread: More potential parallelism, more chance for latency hiding.
  - Tradeoff will depend on GPU resources, typically SM block, thread and register limits.

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 21 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 22 / 22

## Remarks about floating point

- When working on my solution to [last year's HW3](#), Q1,
  - I first wrote:  
 $x = \alpha * x + (1.0 - x);$
  - and the performance was disappointing.
  - After many frustrating attempts to track down the problem, I added one, little  $f$ :  
 $x = \alpha * x + (1.0f - x);$
  - and my code ran 5.5x faster.
- What happened?

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 3 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 4 / 22

## Floating point, doubles, and GPUs

- GPUs are optimized for single-precision floating point arithmetic.
- For the GeForce GTX 550 Ti, double precision arithmetic is way slower than single precision.
- In C, 1.0 is a **double precision** constant, and 1.0f is single precision.
- When I wrote  $x = \alpha * x + (1.0 - x)$ , the compiler generated code that:
  - computes the product  $\alpha * x$ ,
  - both operands are single precision,
  - the computation is done using single precision arithmetic.
  - computes the difference  $1.0 - x$
  - 1.0 is double precision, x is single precision.
  - the computation is done using double precision arithmetic
  - and the result is double precision.
  - computes the product  $\alpha * (1.0 - x)$ .
  - the computation is done using double precision arithmetic
  - and the result is double precision.
- When I wrote  $x = \alpha * x + (1.0f - x)$ , everything stays in single-precision, and it's much faster.

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 7 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 8 / 22

## Memory System Parallelism 1: Get Lots of Data

Memory is addressed per byte, but you retrieve a bunch of (sequential) bytes at once.

- GDDR5 DRAM: 32-bit bus per chip and transfers are in 16 word bursts (so 64 bytes per access per chip).
- GPU global memory (from GDDR DRAMs): Accessed by 32-, 64- or 128-byte transactions.
- Transactions must be "naturally" aligned: First address must be a multiple of the transaction size.
- CC 2.0: L1 cache (1 per SM) serviced by 128-byte transactions, L2 cache (shared by SMs) by 32-byte transactions.
- CC 6.6: Same as 2.0, but L1 cache rules are complicated.
- GPU shared memory (on-chip SRAM): Access in 32-bit words.

Amortize addressing overhead and thereby increase bandwidth.

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 11 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 12 / 22

## Implications for Shared Memory

See [CUDA Toolkit Documentation C Programming Guide Figure 17 and Figure 18](#).

- Consider shared memory address bits:
  - 48KB / thread block requires 16 bits to address.
  - Bottom two bits specify the width in a 32-bit word of data.
  - Next five bits specify which of 32 banks.
  - Top nine bits specify which word within the bank.
- Key takeaway: If two threads in a warp access a memory location in the same bank (same middle five bits of address):
  - If threads access the same location (same top nine bits), then broadcast (on read) or one value wins (on write).
  - If threads access different location, access is serialized (slower but still correct).

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 11 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 12 / 22

## Why all these numbers?

- When designing a new generation of GPUs, the GPU architects run lots of simulations to estimate the performance for various choices of the architectural parameters.
- For example, if more warps are allowed in the scheduling pool
  - The SM will have useful instructions to dispatch more often → better performance.
  - **BUT** the on-chip circuitry to hold and manage the scheduling pool will be larger.
  - This means instruction scheduling will be slower → a longer clock period.
  - Instruction scheduling will use more power → a longer clock period, or fewer SMs, or more expensive chip cooling.
  - The real-state on the chip could have been used for something else. Is this the best use of that area?
  - Note that CUDA 5 made the increase to 64 warps/SM.
- Architects explore these trade-offs to optimize performance for graphics applications, the main source of revenue.
- Architects are also risk-adverse: make the chip as much like the last one that worked as you can.
- These hard-wired constraints have a large impact on program performance.

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 15 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 16 / 22

## Occupancy with CUDA 2.1

- Different GPUs at level CUDA 2.1 have differing numbers of SMs.
- But the SMs all look the same, even for different GPUs.

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 17 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 18 / 22

## SMs, blocks, and threads – the plot

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 19 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 20 / 22

## SMs, blocks, and threads

- Each SM has 32K registers – that's 1K registers per SP.
- This is another constraint:
  - nbks  $\leq \frac{1024}{\text{registersPerThread}}$

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 21 / 22

Greenstreet & Mitchell CUDA: Performance Considerations CPSC 418 – March 17, 2017 22 / 22

## Hitting the register constraint

What if each thread uses 22 registers?

- $22 * 48 = 1056 > 1024 \rightarrow$  can't run 48 warps.
- $\lfloor \frac{1024}{22} \rfloor = \lfloor 46.54 \rfloor = 46$ .
- Can we run 46 warps?
  - One block with 46 warps would have  $46 * 32 = 1472 > 1024$  threads. Not allowed.
  - Two block with 23 warps each would have 736 threads. That should work.
  - But, the plot with the occupancy calculator only shows warp counts that are multiples of 8.
  - Have I overlooked another architectural constraint?
  - probably
- Let's assume that with 23 registers per thread, the SM can run at most 40 warps simultaneously.
  - Then either each thread must have enough instruction-level parallelism to keep the SPs busy.
  - Or, we'll see a drop in performance.

## CUDA: Performance Considerations

Mark Greenstreet & Ian M. Mitchell

CPSC 418 – March 17, 2017

## Thread Divergence

### Thread Divergence

### Floating Point Foibles

### Memory Accesses

### Occupancy

### Granularity

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet & Ian M. Mitchell

and are made available under the terms of the Creative Commons Attribution 4.0 International license

<http://creativecommons.org/licenses/by/4.0/>

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2` uses 17 registers per thread.
- kernel `sh_m1.m1` uses 14 registers per thread.
- both kernels use 4024 bytes of shared memory per block.
- neither kernel spills registers to global memory (good).

© Translation:

- kernel `sh_m1.m2</code`

## Hitting the register constraint

What if each thread uses 22 registers?

- $22 \times 48 = 1056 > 1024 \rightarrow$  can't run 48 warps.
- $\lfloor \frac{1056}{22} \rfloor = \lfloor 46.54 \rfloor = 46$ .
- Can we run 46 warps?
  - One block with 46 warps would have  $46 \times 32 = 1472 > 1024$  threads. Not allowed.
  - Two block with 23 warps each would have each 736 threads. That should work.
  - But, the plot with the occupancy calculator only shows warp counts that are multiples of 8.
  - Have I overlooked another architectural constraint?
  - \* probably
- Let's assume that with 23 registers per thread, the SM can run at most 40 warps simultaneously.
  - Then either each thread must have enough instruction-level parallelism to keep the SP pipelines.
  - Or, we'll see a drop in performance.

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – March 17, 2017 19 / 22

## How many registers does my thread use?

- use the `--ptxas-options -v` option for nvcc
 

```
nvcc --ptxas-options -v -c examples.cu
ptxas info : Compiling entry function '_Z25compileEntryFunction_v' for 'sm_20'
ptxas info : Function properties for _Z25compileEntryFunction_v
0 bytes stack frame, 0 bytes spill stores, 0 bytes spill loads
ptxas info : Compiling entry function '_Z25compileEntryFunction_v' for 'sm_30'
ptxas info : Function properties for _Z25compileEntryFunction_v
0 bytes stack frame, 0 bytes spill stores, 0 bytes spill loads
ptxas info : Used 14 registers, 4096 bytes smem, 56 bytes cmem[0]
```
- Translation:
  - Kernel `shmem2` uses 17 registers per thread.
  - Kernel `shmem1` uses 14 registers per thread.
  - both kernels use 4024 bytes of shared memory per block.
  - neither kernel spills registers to global memory (good).

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – March 17, 2017 20 / 22

## Granularity

How much work should a kernel do?

- Do more work within a kernel: Launching each kernel takes time.
- Do less work within a kernel: New kernels allow for changes in block and grid size, and ensure synchronization between threads even in different blocks.
- Either way: Minimize movement of data to and from the host.
- How much work should a thread do?
- Do more work in a single thread: Fewer chances for memory collisions, easier synchronization, less register contention.
- Do less work in a single thread: More potential parallelism, more chance for latency hiding.
- Tradeoff will depend on GPU resources, typically SM block, thread and register limits.

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – March 17, 2017 21 / 22

## Bigger Kernels

```
global_myKernel(...){
 do something
}
```

Unless `do something` is big, kernel launch takes most of the time.

- We can launch a big-grid
  - If we have a huge number of array elements that each need a small amount of work, this can be a good idea.
  - BUT** we're likely to create a memory-bound problem.
- Or, we can make each thread do many somethings.
 

```
global_myKernel(int m,...){
 for(int i = 0; i < m; i++)
 do something
}
```

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 17, 2017 16 / 18

## Loop Limitations

- It takes two or three instructions per loop iteration to manage the loop:
  - One to update the loop index
  - One or two to check the loop bounds and branch.
  - If `do something` is only three or four instructions, then 40-50% of the execution time is for loop management.
- If each iteration of `do something` depends on the previous one
  - Then the long latency of the SP pipelines can limit performance.
  - Even if we have 48 warps running.

## Preview

|                                                              |
|--------------------------------------------------------------|
| March 20: Matrix multiplication, Part 1                      |
| March 22: Matrix multiplication, Part 2                      |
| March 24: Complete CUDA                                      |
| <b>March 27 – April 3: this may change</b>                   |
| March 27: Using Parallel Libraries                           |
| March 29 – April 3: Verification of and Parallel Programs    |
| April 5: Party, 50 <sup>th</sup> Anniversary of Amdahl's Law |

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 17, 2017 17 / 18

## Loop Unrolling

- Have each loop iteration perform multiple copies of the loop body
 

```
__global__ myKernel(int m, ...){
 for(int i = 0; i < m; i += 4){
 do something
 do something_2
 do something_3
 do something_4
 }
}
```
- More "real work" for each time the loop management code is executed.
- Need to make sure that `m` is a multiple of four, or handle end-cases separately.
- Often, we need more registers.

This example is from [last year's HW3](#), Q1.

## Matrix Multiplication – Algorithms

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 17, 2017 18 / 18

### Mark Greenstreet

CpSc 418 – Mar. 22, 2017

#### Outline:

- Sequential Matrix Multiplication
- Parallel Implementations, Performance, and Trade-Offs.

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license (<http://creativecommons.org/licenses/by/4.0/>)

## Sequential Matrix Multiplication

```
mult(A, B) ->
 BT = transpose(B),
 lists:map(
 fun(RowOfB) ->
 lists:map(
 fun(ColOfB) ->
 dot_prod(RowOfA, ColOfB)
 end, BT
 end, A)
 dot_prod(V1, V2) ->
 lists:foldl(
 fun({X,Y},Sum) -> Sum + X*Y end,
 0, lists:zip(V1, V2)).
```

Next, we'll see list comprehensions to get a more succinct version.

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 22 / 22

## Matrix Multiplication, with comprehensions

```
mult(A, B) ->
 BT = transpose(B),
 [lists:map(
 fun(ColOfB) ->
 dot_prod(RowOfA, ColOfB)
 end, BT
 end, A)
transpose([]) -> [] ; % special case for empty matrices
transpose([_|_]) -> [_| transpose([_|_])] ; % bottom of recursion, the columns are empty
transpose([_|_]) -> [| [H | T] <- [| H] % create a row from the first column of M
 | transpose([| [H | T] <- [| H] % now, transpose what's left
].].
```

## Tiling Matrices

### An Example

- Let `A`, `B`, and `C = AB` be  $16 \times 16$  matrices.
- Let `A1 = A[1..4, 1..16]`, i.e. the first four rows of `A`.
  - In our Erlang representation, `[A1, ...] = lists:split(4, A)`.
- Let `A2 = A[5..8, 1..16]`; `A3 = A[9..12, 1..16]`;
 `A4 = A[13..16, 1..16]`; and likewise for `C1, C2, C3`, and `C4`.
- Big important fact:

$$\begin{aligned} C1 &= A1 \cdot B & C2 &= A2 \cdot B \\ C3 &= A3 \cdot B & C4 &= A4 \cdot B \end{aligned}$$

### In sequential Erlang:

$$[C1, C2, C3, C4] = [mult(AA, B) || AA <- A]$$

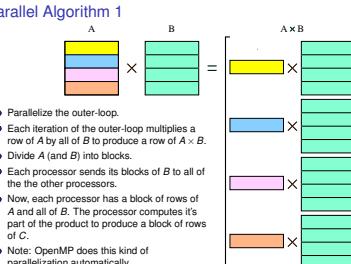
- To make it parallel, we compute each of the `Ci = Ai · B` with a separate process.

## Performance of Parallel Algorithm 1 – Measured

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 8 / 20

### Parallel Algorithm 1

- Parallelize the outer-loop.
- Each iteration of the outer-loop multiplies a row of `A` by all of `B` to produce a row of `A · B`.
- Divide `A` and `B` into blocks.
- Each processor sends its blocks of `B` to all of the other processors.
- Now, each processor has a block of rows of `A` and all of `B`. The processor computes its part of the product to produce a block of rows of `C`.
- Note: OpenMP does this kind of parallelization automatically.

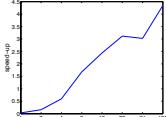


## Performance of Parallel Algorithm 2

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 12 / 20

### Using Memory more Efficiently

- Main idea: each process works on one "slab" of `B` at a time.
- $C[i,j] = \sum_{k=1}^N A[i,k] B[k,j]$ , a dot-product
- $= (\sum_{k=1}^{N/4} A[i,k] B[k,j]) + (\sum_{k=N/4+1}^{N/2} A[i,k] B[k,j]) + (\sum_{k=N/2+1}^{3N/4} A[i,k] B[k,j]) + (\sum_{k=3N/4+1}^{4N} A[i,k] B[k,j])$
- Each process does each of its four summations when it holds the corresponding slab of `B`.
  - Each holds one slab of `A` for the whole computation.
  - Each process only needs to hold one slab of `B` at a time.
- The algorithm generalizes to having any number of slabs for `A` and `B` in the obvious way.
  - Should be "obvious" if I've explained this clearly.
  - If it isn't obvious, that's my bad – please ask a question.

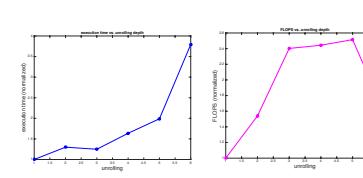


## Bad performance, pass it on

- CPU operations: Same as for parallel algorithm 1: total time:  $\mathcal{O}(N^2/P)$ .
- Communication: Same as for parallel algorithm 1:  $\mathcal{O}(N^2 + P)$ .
  - With algorithm 1, each processor sent the same message to  $P - 1$  different processors.
  - With algorithm 2, for each processor, there is one destination to which it sends  $P - 1$  different messages.
  - Thus, algorithm 2 can work efficiently with simpler interconnect networks.
- Memory: Each process needs  $\mathcal{O}(N^2/P)$  storage for its block of `A`, its current block of `B`, and its block of the result.
- Note: each processor might hold onto its original block of `B` so we still have the blocks of `B` available at the expected processors for future operations.
- Do the memory savings matter?

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 16 / 20

## Unrolling – the plots



This example is from [last year's HW3](#), Q1.

## Objectives

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 19 / 19

## Where's A?

- Communication between the CPU and GPU
  - Kernel launch overhead
  - Transferring data between CPU memory and GPU memory
    - Is this solved with more recent GPUs that can access the GPU memory directly?
    - Not really, the data still needs to be transferred.
    - And it's one more memory level for the programmer to keep track of.
- Communication between blocks
  - Write global memory and end the kernel
  - Launch a new kernel and read the global memory.
  - The same strategy applies if the shape for the required grid changes between phases of a larger computation.
- Communication between warps in a block
  - `synchthreads`
- AND
  - There's a built-in energy cost of the big register file.
  - Trade-offs of energy, latency, and parallelism. large numbers of threads.

## Matrix representation in Erlang

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 2 / 22

I'll represent a matrix as a list of lists.

For example, the matrix

$$\begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 4 & 9 & 16 \\ 1 & 8 & 27 & 64 \end{bmatrix}$$

is represented by the Erlang nested-list:

$$\begin{bmatrix} [1, 2, 3, 4] \\ [1, 4, 9, 16] \\ [1, 8, 27, 64] \end{bmatrix}$$

The empty matrix is `[]`.

- This means my representation can't distinguish between a  $2 \times 0$  matrix, a  $0 \times 4$  matrix, and a  $0 \times 0$  matrix.
- That's OK. This package is to show some simple examples.
- I'm not claiming it's for advanced scientific computing.

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 3 / 20

## Performance – Modeled

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 4 / 20

### Performance – Modeled

- Really simple, operation counts:
  - Multiplications:  $n_{cols} \cdot n_{rows} \cdot n_{cols\_b} \cdot n_{cols\_a}$ .
  - Additions:  $n_{rows} \cdot n_{cols\_b} \cdot n_{cols\_a} - 1$ .
  - Memory reads:  $2 \cdot N^2$  multiplications.
  - Memory writes:  $n_{rows} \cdot n_{cols\_a} \cdot n_{cols\_b}$ .
  - Time is  $\mathcal{O}(n_{rows} \cdot n_{cols\_a} \cdot n_{cols\_b} \cdot n)$ . If both matrices are  $N \times N$ , then it's  $\mathcal{O}(N^3)$ .
- But, memory access can be terrible.
  - For example, let matrices `a` and `b` be  $1000 \times 1000$ .
  - Assume a processor with a 4M L2-cache (final cache), 32 byte-cache lines, and a 200 cycle stall for main memory accesses.
  - Observe that a row of `matrix a` and a column of `b` fit in the cache. (a total of ~40K bytes).
  - But, all of `a` does not fit in the cache (that's 8 Mbytes).
  - So, on every fourth pass through the inner loop, `every` read from `b` is a cache miss!
  - Cache miss dominates everything else.
- This is why there are carefully tuned numerical libraries.

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 5 / 20

## Performance – Measured

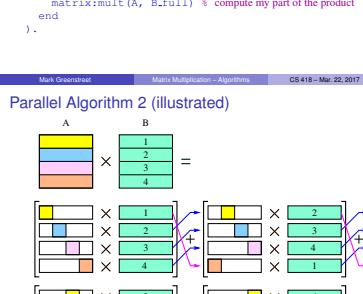
Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 6 / 20

### Performance – Measured

- CPU operations: same total number of multiplies and adds, but distributed around  $P$  processors. Total time:  $\mathcal{O}(N^2/P)$ .
- Communication: Each processor sends (and receives)  $P - 1$  messages of  $N^2/P$ . If time to send a message is  $t_0 + t_1 \cdot M$  where  $M$  is the size of the message, then the communication time is
 
$$(P - 1) \left( t_0 + \frac{N^2}{P} \right) = \mathcal{O}(N^2 + \lambda P)$$
- Note: I'm assuming  $t_0$  corresponds to  $\lambda$ , and that  $t_1$  is roughly the same as a the time for "typical" sequential operations..
- Memory: Each process needs  $\mathcal{O}(N^2/P)$  storage for its block of `A` and the result. It also needs  $\mathcal{O}(N^2)$  to hold all of `B`.
  - The simple algorithm divides the computation across all processors, but it doesn't make good use of their combined memory.

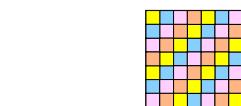
Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 7 / 20

## Parallel Algorithm 2 (illustrated)



## Tiling in Real-Life

- Each processor first computes what it can with its rows from `A` and `B`.
  - It can only use  $N/P$  of its columns from `B`.
  - It uses its entire block from `B`.
  - We've now computed one of `P` matrices, where the sum of all of these matrices is the matrix `C`.
- We view the processors as being arranged in a ring.
  - Each processor forwards its block of `B` to the next processor in the ring.
  - Each processor computes a new partial product of `AB` and adds it to what it had from the previous step.
  - This process continues until every block of `B` has been used by every processor.
- Why? If there's time, I'll explain in class.



## Summary

- Matrix multiplication is well-suited for a parallel implementation.
- Need to consider communication costs.
- In the previous algorithms, compute time grows as  $N^2/P$ , while communication time goes as  $(N^2 + P)$ .
- Thus, if  $N$  is big enough, computation time will dominate communication time.
- Connection of theory with actual run time is pretty good:
  - But the matrices have to be big enough to amortize the communication costs.

Greenstreet & Mitchell CUDA: Performance Considerations CS418 – Mar. 22, 2017 19 / 20

|                                                                     |
|---------------------------------------------------------------------|
| <b>March 24:</b> Matrix Multiplication in CUDA                      |
| Homework: HW due at 11:59pm<br>HW5 goes out                         |
| <b>March 27:</b> Using Parallel Libraries                           |
| <b>March 29:</b> Introduction to Model Checking                     |
| Reading: TBA                                                        |
| <b>March 31:</b> The PReach Model Checker                           |
| Reading: Industrial Strength...Model Checking                       |
| <b>April 3:</b> Distributed Termination Detection                   |
| <b>April 5:</b> Party: 50 <sup>th</sup> Anniversary of Amdahl's Law |

## CUDA: Matrix Multiplication

Mark Greenstreet

CpSc 418 – Mar. 24, 2017

- A Brute Force Implementation
- Tiling

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license  
<http://creativecommons.org/licenses/by/4.0/>

Mark Greenstreet Matrix Multiplication - Algorithms CS 418 – Mar. 22, 2017 20 / 20

Mark Greenstreet CUDA: Matrix Multiplication CS 418 – Mar. 24, 2017 1 / 12

## Brute-force performance

- Not very good – each loop iteration performs
  - Two global memory reads.
  - One fused floating-point add.
  - Four or five integer operations.
- Global memory is slow
  - Long access times.
  - Bandwidth shared by all the SPs.
- This implementation has a low CGMA
  - CGMA = Compute to Global Memory Access ratio  $\approx 1/2$ .
- Performance should be:
  - $O(n^3)$  computation,  $O(N^2)$
  - wall-clock:  $\sim O(N^3)$  with  $\alpha$  determined mainly by global memory bandwidth.
  - measured:  $T(1024) \approx 0.0986s$ ;  $T(2048) \approx 0.797s$ ;  $T(3072) \approx 2.7s$ ;  $T(4096) \approx 6.3s$ .
  - $N^3 / T(N) \approx 11 / ns$  – i.e. about  $20 \times 10^9$  multiply-adds per second. Well below GPU peak floating point capacity. Demonstrates global memory bandwidth bottleneck (with a little help from the on-chip caches).

## Tiles vs. Slabs



## Tiles vs. Slabs



- Matrix multiplication: each processor (color) has tiles at (i,j) and (j,i)
  - Can compute all products for the main diagonal, and stripes at spacings of  $P$ .
- Use a reduce to combine results to get the main diagonal and the stripes.
- Rotate  $B$  one block to the left, and compute the next set of stripes.
- After  $P$  rotations, the computation is done.
- Same amount of work (and communication) as the improved slab method from Wednesday.
- Other algorithms such as LU-Decomposition
  - Rows and columns are eliminated from the left and the top.
  - Tiles provide better load balancing.

Mark Greenstreet CUDA: Matrix Multiplication CS 418 – Mar. 24, 2017 4 / 12

Mark Greenstreet CUDA: Matrix Multiplication CS 418 – Mar. 24, 2017 4 / 12

Mark Greenstreet CUDA: Matrix Multiplication CS 418 – Mar. 24, 2017 3 / 12

## Tiling the computation

- Divide each matrix into  $m \times m$  tiles.
- For simplicity, we'll assume that  $n$  is a multiple of  $m$ .
- Each block computes a tile of the product matrix.
- Computing a  $m \times m$  tile involves computing  $n/m$  products of  $m \times m$  tiles and summing up the results.

## Performance issues for mmult2

## The "checklist"

- Are global memory accesses coalesced?
- What is the CGMA?
- Do we have shared memory access conflicts?
- What is the warp-scheduler occupancy?
- How many registers per thread?
- How many threads per block?
- How much shared memory per block?
- How much "other stuff" does each thread perform for each floating point operation?

## Linear Algebra Libraries and CUDA

Mark Greenstreet &amp; Ian M. Mitchell

CPSC 418 – March 27, 2017

## Why?

- BLAS
- Using BLAS (in general)
- Using BLAS (on CUDA GPUs)
- Other numerical libraries

Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet & Ian M. Mitchell  
<http://creativecommons.org/licenses/by/4.0/>

## Levels of BLAS

BLAS specification consists of operations at one of three "levels":

- BLAS-1: Vector-vector operations (scalar vector product, vector sum, dot product, etc.)
  - [Lawson et al., 1979].
  - Performs  $O(n)$  operations on  $O(n)$  data.
- BLAS-2: Matrix-vector operations (matrix-vector product, triangular solves)
  - [Dongarra et al., 1988].
  - Performs  $O(r^2)$  operations on  $O(r^2)$  data.
- BLAS-3: Matrix-matrix operations (matrix-matrix product, triangular solves with multiple right-hand sides)
  - [Dongarra et al., 1990].
  - Performs  $O(r^3)$  operations on  $O(r^3)$  data.

## Deciphering BLAS Function Arguments (part 1)

Consider matrix product  $C = \alpha A^T B^H + \beta C$  implemented by

```
blas_sgemm(enum blas_orderType layout,
 enum blas_transType transa,
 enum blas_transType transb,
 int m, int n, int k,
 float alpha,
 float *a, int lda,
 float *b, int ldb,
 float beta,
 float *c, int ldc)
```

- layout specifies either column-major or row-major.
- transa specifies whether to use  $A$ ,  $A^T$  or  $A^H$ .
- Same for transb and B.
- m, n, k specify matrix sizes:  $A$  is  $m \times k$ ,  $B$  is  $k \times n$ ,  $C$  is  $m \times n$ .
- alpha and beta specify scalar multipliers.
- Some implementations may require pass by reference.

Greenstreet &amp; Mitchell Linear Algebra Libraries and CUDA CPSC 418 – March 27, 2017 6 / 15

## CUDA: Matrix Multiplication

Mark Greenstreet

CpSc 418 – Mar. 24, 2017

- A Brute Force Implementation
- Tiling

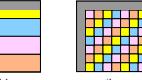
Unless otherwise noted or cited, these slides are copyright 2017 by Mark Greenstreet and are made available under the terms of the Creative Commons Attribution 4.0 International license  
<http://creativecommons.org/licenses/by/4.0/>

Mark Greenstreet Matrix Multiplication - Algorithms CS 418 – Mar. 22, 2017 20 / 20

Mark Greenstreet CUDA: Matrix Multiplication CS 418 – Mar. 24, 2017 1 / 12

Mark Greenstreet CUDA: Matrix Multiplication CS 418 – Mar. 24, 2017 2 / 12

## Tiles vs. Slabs



- Matrix multiplication: each processor (color) has tiles at (i,j) and (j,i)
  - Can compute all products for the main diagonal, and stripes at spacings of  $P$ .
- Use a reduce to combine results to get the main diagonal and the stripes.
- Rotate  $B$  one block to the left, and compute the next set of stripes.
- After  $P$  rotations, the computation is done.
- Same amount of work (and communication) as the improved slab method from Wednesday.
- Other algorithms such as LU-Decomposition
  - Rows and columns are eliminated from the left and the top.
  - Tiles provide better load balancing.

Mark Greenstreet CUDA: Matrix Multiplication CS 418 – Mar. 24, 2017 4 / 12

Mark Greenstreet CUDA: Matrix Multiplication CS 418 – Mar. 24, 2017 4 / 12

Mark Greenstreet CUDA: Matrix Multiplication CS 418 – Mar. 24, 2017 4 / 12

## A Tiled Kernel (step 1)

## A Tiled Kernel (step 1)

```
#define TILEWIDTH 16
__global__ mmult2(float *a, float *b, float *c, int n) {
 float *a_row = a + (blockDim.y*blockIdx.y + threadIdx.y)*n;
 float *b_col = b + (blockDim.x*blockIdx.x + threadIdx.x)*n;
 float sum = 0.0;
 for(int k1 = 0; k1 < gridDim.x; k1++) { % each tile product
 for(int k2 = 0; k2 < blockDim.x; k2++) { % within each tile
 k = k1*blockDim.y + k2;
 sum += a_row[k] * b_col[k];
 }
 }
 c[(blockDim.y*blockIdx.y + threadIdx.y)*n +
 (blockDim.x*blockIdx.x + threadIdx.x)] = sum;
}
```

## Launching the kernel:

```
int nbk = n/TILEWIDTH;
dim3 blks(nbk, nbk, 1);
dim3 thrd3s(TILEWIDTH, TILEWIDTH, 1);
matrixMult<>blks,thrd3s>>(a, b, c, n);
```

Other numerical applications:

- LU-decomposition and other factoring algorithms.
- Matrix transpose.
- Finite element methods.
- Many, many more.

A non-numerical example: revsort

- To sort  $N^2$  values, arrange them as a  $N \times N$  array.
 

```
repeat logN times {
 sort even numbered rows left-to-right.
 sort odd numbered rows right to left.
 sort columns top-to-bottom.
}
```
- We can get coalesced accesses for the rows, but not the columns.
  - Cooperative loading can help here – e.g. use a transpose.

## Performance of mmult2

## Performance of

## Notes on cuBLAS

- Always uses column-major ordering
  - So be careful with data layout.
- Always uses 1-based indexing.
  - Usually irrelevant since you do not index into arrays.
- All cuBLAS code is called from the host.
  - You do not write any kernel code.
  - You do not have to worry about grids, blocks, shared memory, ...
- Need to link against cuBLAS library.
  - Check that environment variable `LD_LIBRARY_PATH` includes CUDA library directory.
  - `(/usr/local/lib/pkg/cudatoolkit/lib64 on linXX machines)`
  - Add `-lcublas` to compile command.

## Efficiency of cuBLAS

Matrix product example from [2017-03-24 lecture](#) (in seconds):

|                                      | 1024  | 2048  | 3072  | 4096  |
|--------------------------------------|-------|-------|-------|-------|
| Brute force <code>mnmmt</code>       | 0.079 | 0.648 | 2.190 | 5.152 |
| Tiled <code>mnmmt<sup>2</sup></code> | 0.027 | 0.208 | 0.724 | 1.690 |
| <code>cublas.sgemm</code>            | 0.007 | 0.052 | 0.176 | 0.421 |

- Brute force `mnmmt` achieves  $\sim 13$  GFLOPS.
- `cublas.sgemm` achieves  $\sim 160$  GFLOPS.
- (Nolan's tiled implementation `mnmmt2` is buggy.)

Greenstreet & Mitchell

Linear Algebra Libraries and CUDA

CPSC 418 – March 27, 2017

13 / 15

Greenstreet & Mitchell

Linear Algebra Libraries and CUDA

CPSC 418 – March 27, 2017

14 / 15

Greenstreet & Mitchell

Linear Algebra Libraries and CUDA

CPSC 418 – March 27, 2017

15 / 15

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

1 / 13

## Model-Checking: Motivation

- What is “model checking”?
  - Construct a “model” for a piece of hardware or software – typically a finite-state machine.
  - Give a precise, mathematical definition of properties that the design is supposed to have.
  - Show that that model satisfies the specification.
    - For example, find all reachable states of the model.
    - Show that every reachable state satisfies a desired property – for example, mutual exclusion.
- Why use model checking?
  - Find bugs.
  - Hardware bugs are very expensive.
  - Software bugs are very common, but
    - Finding bugs in concurrent software is **hard**.
    - The challenges of finding bugs motivates using more systematic approaches.
- A simple example: Dekker’s Mutual Exclusion algorithm

## Dekker’s Algorithm

Problem statement: ensure that at most one thread is in its critical section at any given time.

```
thread 0: thread 1:
PC0= 0: while(true) { PC0= 0: while(true) {
PC0= 1: non-critical code PC0= 1: non-critical code
PC0= 2: flag[0] = true; PC0= 2: flag[1] = true;
PC0= 3: while(flag[0]) { PC0= 3: while(flag[0]) {
PC0= 4: if(turn!= 0) { PC0= 4: if(turn != 1) {
PC0= 5: flag[0] = false; PC0= 5: flag[1] = false;
PC0= 6: } PC0= 6: while(turn != 0);
PC0= 7: turn = 1; PC0= 7: turn = 0;
PC0= 8: flag[0] = true; PC0= 8: flag[1] = true;
PC0= 9: } PC0= 9: }
PC0=10: critical section PC0=10: critical section
PC0=11: turn = 1; PC0=11: turn = 0;
PC0=12: flag[0] = false; PC0=12: flag[1] = false;
PC0=13: } PC0=13: }
```

See [http://en.wikipedia.org/wiki/Dekker's\\_algorithm](http://en.wikipedia.org/wiki/Dekker's_algorithm).

## Model Checking Dekker’s algorithm

- Represent each state with 9-bits:
  - three for the location of each process (6 locations)
  - three for flag and turn
  - I’ll show a simple python version that uses python tuples
- Pseudo-code:

```
initialState = (1,1,0,0,0); // (loc0, loc1, flag0, flag1, turn)
workList = queue(); // initially empty
knownStates = set(); // initially empty
workList.insert(initialState);
while len(workList) > 0:
 s = workList.removeNext();
 if s[4] == 0: // for mutual exclusion:
 for s' in nextStates(s):
 if s' not in knownStates:
 add s' to worklist and knownStates;
```
- Model-checking finds 48 reachable states for Dekker’s algorithm and verifies mutual exclusion.

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

2 / 13

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

3 / 13

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

4 / 15

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

5 / 13

## murφ: a guarded command language

## murφ: execution model

- In murφ a guarded command is called a rule and is written:

`rule guard => action`

► When `guard` is satisfied, `action` may be performed.

► Example: rule (`(loc[0] == 3) and flag[1]`)  $\rightarrow$  `loc[0] := 4`

- Rules may be quantified using the `Ruleset` construction:

```
Process: scalarSet2();
ruleset: Process do
 (loc[i] == 3) and flag[i] => loc[i] := 4
 end
```

- Model-checking finds 48 reachable states for Dekker’s algorithm and verifies mutual exclusion.

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

6 / 13

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

7 / 13

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

8 / 15

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

9 / 13

## Review

## Industrial Strength, Parallel Model Checking

Mark Greenstreet

CpSC 418 – Mar. 31, 2017

I'll add some review questions.

## Industrial Strength, Parallel Model Checking

## Industrial Strength, Parallel Model Checking

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

10 / 13

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

11 / 13

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

12 / 13

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

13 / 13

## Model Checking: the algorithm

## Model Checking: the algorithm

Mark Greenstreet

Model Checking

CS 418 – Mar. 29, 2017

14 / 13

Mark Greenstreet

Industrial Strength, Parallel Model Checking

CS 418 – Mar. 31, 2017

1 / 17

Mark Greenstreet

Industrial Strength, Parallel Model Checking

CS 418 – Mar. 31, 2017

2 / 17

Mark Greenstreet

Industrial Strength, Parallel Model Checking

CS 418 – Mar. 31, 2017

2 / 17

## Model Checking: the algorithm

## Model Checking: the algorithm

Mark Greenstreet

Model Checking

CS 418 – Mar. 31, 2017

3 / 17

Mark Greenstreet

Industrial Strength, Parallel Model Checking

CS 418 – Mar. 31, 2017

3 / 17

Mark Greenstreet

Industrial Strength, Parallel Model Checking

CS 418 – Mar. 31, 2017

3 / 17

Mark Greenstreet

Industrial Strength, Parallel Model Checking

CS 418 – Mar. 31, 2017

3 / 17

## How can we make this algorithm parallel?

How can we make this algorithm parallel?

- Where does the code spend most of the time? `next_states(s)`

► Where does the code spend most of the time? `next_states(s)`

- What are the dependencies?

► As written, the code is very sequential.

► BUT, the correctness of the algorithm does not depend on the order in which states are removed from `workList`.

- What uses most of the memory? `knownStates` and `workList`

How can we make this algorithm parallel?

- Where does the code spend most of the time? `next_states(s)`

► Where are the dependencies?

► As written, the code is very sequential.

► BUT, the correctness of the algorithm does not depend on the order in which states are removed from `workList`.

► What uses most of the memory? `knownStates` and `workList`

## Model Checking: the algorithm

## Model Checking: the algorithm

Mark Greenstreet

Model Checking

CS 418 – Mar. 31, 2017

3 / 17

Mark Greenstreet

Industrial Strength, Parallel Model Checking

CS 418 – Mar. 31, 2017

3 / 17

Mark Greenstreet

Industrial Strength, Parallel Model Checking

CS 418 – Mar. 31, 2017

4 / 17

Mark Greenstreet

Industrial Strength, Parallel Model Checking

CS 418 – Mar. 31, 2017

5 / 17

## Notes on cuBLAS

## Efficiency of cuBLAS

Matrix product example from [2017-03-24 lecture](#) (in seconds):

|                                      | 1024  | 2048  | 3072  | 4096  |
|--------------------------------------|-------|-------|-------|-------|
| Brute force <code>mnmmt</code>       | 0.079 | 0.648 | 2.190 | 5.152 |
| Tiled <code>mnmmt<sup>2</sup></code> | 0.027 | 0.208 | 0.724 | 1.690 |
| <code>cublas.sgemm</code>            | 0.007 | 0.052 | 0.176 | 0.421 |

- Brute force `mnmmt` achieves  $\sim 13$  GFLOPS.
- `cublas.sgemm` achieves  $\sim 160$  GFLOPS.
- (Nolan’s tiled implementation `mnmmt2` is buggy.)

## Model-Checking: Motivation

## Dekker’s Algorithm

Problem statement: ensure that at most one thread is in its critical section at any given time.

thread 0: thread 1:

```
PC0= 0: while(true) { PC0= 0: while(true) {
PC0= 1: non-critical code PC0= 1: non-critical code
PC0= 2: flag[0] = true; PC0= 2: flag[1] = true;
PC0= 3: while(flag[0]) { PC0= 3: while(flag[0]) {
PC0= 4: if(turn!= 0) { PC0= 4: if(turn != 1) {
PC0= 5: flag[0] = false; PC0= 5: flag[1] = false;
PC0= 6: } PC0= 6: while(turn != 0);
PC0= 7: turn = 1; PC0= 7: turn = 0;
PC0= 8: flag[0] = true; PC0= 8: flag[1] = true;
PC0= 9: } PC0= 9: }
PC0=10: critical section PC0=10: critical section
PC0=11: turn = 1; PC0=11: turn = 0;
PC0=12: flag[0] = false; PC0=12: flag[1] = false;
PC0=13: } PC0=13: }
```

See [http://en.wikipedia.org/wiki/Dekker's\\_algorithm](http://en.wikipedia.org/wiki/Dekker's_algorithm).

## Model Checking Dekker’s algorithm

## A Brief History of Model Checking

Proposed by Clarke and Emerson (1981) and independently by Sifakis (1982).

► They shared the 2007 Turing Award.

► Their approach was essentially the one described above.

Symbolic methods introduced by McMillan (1987) using binary-decision diagrams, a DAG representation of boolean formulas.

Widespread adaptation of model-checking for hardware design took place in the 1990s and continues today.

The murφ model checker is a landmark in this work.

Model-checking of software is now gaining industrial acceptance

- Based on “predicate abstraction” methods of Clarke and Grumberg, and independently Ball.

- Enabled in aerospace, SAT solvers and interpolation-based model checking (McMillan).

## Model Checking Dekker’s algorithm

## Is Dekker’s algorithm correct?

Problem statement: ensure that at most one thread is in its critical section at any given time.

thread 0: thread 1:

```
PC0= 0: while(true) { PC0= 0: while(true) {
PC0= 1: non-critical code PC0= 1: non-critical code
PC0= 2: flag[0] = true; PC0= 2: flag[1] = true;
PC0= 3: while(flag[0]) { PC0= 3: while(flag[0]) {
PC0= 4: if(turn!= 0) { PC0= 4: if(turn != 1) {
PC0= 5: flag[0] = false; PC0= 5: flag[1] = false;
PC0= 6: } PC0= 6: while(turn != 0);
PC0= 7: turn = 1; PC0= 7: turn = 0;
PC0= 8: flag[0] = true; PC0= 8: flag[1] = true;
PC0= 9: } PC0= 9: }
PC0=10: critical section PC0=10: critical section
PC0=11: turn = 1; PC0=11: turn = 0;
PC0=12: flag[0] = false; PC0=12: flag[1] = false;
PC0=13: } PC0=13: }
```

See [http://en.wikipedia.org/wiki/Dekker's\\_algorithm](http://en.wikipedia.org/wiki/Dekker's_algorithm).

## Model Checking Dekker’s algorithm

## Is Dekker’s algorithm correct?

Problem statement: ensure that at most one thread is in its critical section at any given time.

thread 0: thread 1:

```
PC0= 0: while(true) { PC0= 0: while(true) {
PC0= 1: non-critical code PC0= 1: non-critical code
PC0= 2: flag[0] = true; PC0= 2: flag[1] = true;
PC0= 3: while(flag[0]) { PC0= 3: while(flag[0]) {
PC0= 4: if(turn!= 0) { PC0= 4: if(turn != 1) {
PC0= 5: flag[0] = false; PC0= 5: flag[1] = false;
PC0= 6:
```

**Parallel Model Checking: Performance Analysis**

- A sequential implementation of the model checking algorithm requires  $\mathcal{O}(SR)$  time, where  $S$  is the number of reachable states, and  $R$  is the number of rules.
- A parallel implementation requires  $\mathcal{O}(SR/P)$  compute time, and sends worst-case  $\mathcal{O}(SR)$  messages.
  - In practice, the average number of successors of each state (i.e. the degree of the state-graph) is relatively small. If we assume this is a small constant, then we get  $\mathcal{O}(S)$  messages.
- Consider the case where each worker process generates  $\sigma$  new successor states,  $\mathcal{O}(\sigma)$ , per second.
  - These are sent to the other processes uniformly at random (if we have a good hash function).
  - Half of these messages cross the bisection of any network.
  - That means we have a bottleneck bandwidth of  $\approx P/2$ .
- For real-life networks, bisection bandwidth grows much more slowly than  $P$ .
  - If we scale this algorithm to a large enough number of processors, network bandwidth will be the limiting constraint.
  - This is a common performance pattern in parallel computing.

Mark Greenstreet Industrial Strength: Parallel Model Checking CS 418 – Mar. 31, 2017 6 / 17

## Batching Messages

- If we send each new state,  $s^*$ , one at a time to its owner process, communication overhead dominates the run time.
- Key lesson:** pay attention to  $\lambda$ .
- The Erlang code maintains a separate buffer for each worker process.
  - Add states that should be sent to that process to the buffer until we have enough.
  - Then send them as a batch.
  - A process that is running low on work sends requests to the other workers to ask them to flush their buffers.
  - These flush requests are bundled with state batches to avoid extra messages.
- Erlang makes the communication architecture simple and easy to extend.

Mark Greenstreet Industrial Strength: Parallel Model Checking CS 418 – Mar. 31, 2017 10 / 17

## Flexibility

- The Erlang code for PRReach is simple.
  - The version described in the paper is about 1000 lines of code.
- This makes PRReach a flexible platform for experiments:
  - Checking response properties: e.g. every request is eventually granted.
  - Exploiting symmetry: there are times we can verify a protocol for two or three nodes and conclude with certainty that it is correct for any number of nodes.
  - And others.
- Applications:
  - Used by Intel architects when exploring protocols for on-chip memory.
  - Used in other companies and universities.
  - It's been run on hundreds of machines with models of hundreds of billions of states.
  - Symbolic methods are faster than PRReach for safety properties (showing that the model never reaches a bad state).
    - PRReach is faster for handling liveness properties: showing that some condition will eventually be satisfied.

Mark Greenstreet Industrial Strength: Parallel Model Checking CS 418 – Mar. 31, 2017 14 / 17

## Review: for today's lecture

- To make a parallel implementation of a computation, we often need to identify a set (or many sets!) of operations that can be performed in parallel. Does this exactly replicate the sequential version, or does it perform an equivalent computation?
- Same questions as above, but for reduce.
- Scan? Sorting? Matrix multiplication?
- Why did slow processes in PRReach tend to become catastrophically slower? What was the solution?
- What is load balancing? Compare the load balancing mechanisms of PRReach and Google's map-reduce.
- Is Erlang suitable for large-scale, high-performance, parallel computing? Why or why not?

Mark Greenstreet Industrial Strength: Parallel Model Checking CS 418 – Mar. 31, 2017 18 / 17

## Parallel Algorithms 2

- Matrix operations
  - matrix multiplication dividing the matrix into b
    - Dividing the matrix into blocks
    - Analysis of the compute and communication costs.
  - BLAS and cuBLAS: use a library when you can!
- Map-Reduce
- Model checking
  - We only had two lectures on the topic and no HW.
  - Won't ask any detailed questions, but if there might be some high-level questions with one sentence answers in the review questions, in which case similar questions could be on the exam.
  - Know what model checking is: verifying properties of a hardware or software design, using a finite-state-machine model, finding the reachable states.
  - The idea of distributing work by hashing values and sending each to its owner process.

Mark Greenstreet Course Summary CS 418 – Apr. 5, 2017 3 / 8

## ... and the things we have left undone

- More paradigms and programming frameworks
  - shared memory: Java threads, pthreads.
  - futures: e.g. Scala
  - MPI and OpenMP (for scientific computing)
  - many big-data, machine-learning, and scientific computing frameworks
- Do a bigger project.
- The good news:**
  - You've got what you need to learn new paradigms, new frameworks, and take on realistic projects.
  - From my experience with research projects that have moved into industry in the past few years, you've got the critical knowledge and skills.
  - Writing industrial-strength, parallel-code with good performance is still more than a homework assignment, but when my students have done it, they've built on the foundation you now have.

Mark Greenstreet Course Summary CS 418 – Apr. 5, 2017 7 / 8

## From Algorithm to Industrial Adaptation

What is needed for real-world verification?

- Lots of memory:
  - Memory and time are both concerns for model checkers, but memory tends to be the more critical concern.
  - A parallel implementation offers the combined memory of a large number of machines.
- Robustness:
  - Simple architecture and re-use stable, well-exercised code.
  - Prevent "overwhelm and crash".
  - Load balancing.
- Flexibility:
  - Solve problems that other tools cannot
  - In particular, liveness properties such as "response".

## Overwhelm and Crash

The dangers of using Erlang for high-performance computing

- The Erlang inbox is a list.
  - Newly received messages are prepended to the list.
  - A receive gets the oldest message that matches a pattern of the receive.
  - This means that the time for receive is linear in the number of pending message.
- This leads to a performance catastrophe
  - If a process gets slightly behind, its inbox will fill a little more than the other processes.
  - This means that a process that falls behind will slow down, and its inbox will fill even more.
  - Eventually, the process crashes.

## Termination

How do we know when we're done?

- Well, times up for this lecture.
- More seriously, in PRReach we need to know when
  - When every worker process has an empty worklist.
  - And there are no messages in flight.
- Both conditions must hold at the same time.
  - This is the topic for Monday's lecture.

## Review: for March 29 lecture

- What is model checking?
- What is mutual exclusion?
- How does the model checker presented in the March 29 slides show that Dekker's algorithm guarantees mutual exclusion.
- Describe the role of the `knownStates` and `workList` data structures in the model checking algorithm.
- What is a guarded command?
- Write a  $\mu\text{-}\varphi$  rule for another statement from Dekker's algorithm.

## Parallel Architectures

- pipelining and instruction level parallelism
- shared memory multiprocessors
  - Know what a cache coherence protocol is.
  - Explain the idea of shared-reader or exclusive writer.
  - Be able to point out that real cache coherence protocols aren't as "consistent" as the simple (e.g. MESI) model from class.
- message passing architectures
  - Rings, tori, hypercubes
  - Latency, bisection width.
- data parallel architectures
  - SIMD (and SIMD)
  - instruction execution: why so many threads
  - GPU memory hierarchy

Mark Greenstreet Course Summary CS 418 – Apr. 5, 2017 4 / 8

## ... and the things we have left undone

- More paradigms and programming frameworks
  - shared memory: Java threads, pthreads.
  - futures: e.g. Scala
  - MPI and OpenMP (for scientific computing)
  - many big-data, machine-learning, and scientific computing frameworks
- Do a bigger project.
- The good news:**
  - You've got what you need to learn new paradigms, new frameworks, and take on realistic projects.
  - From my experience with research projects that have moved into industry in the past few years, you've got the critical knowledge and skills.
  - Writing industrial-strength, parallel-code with good performance is still more than a homework assignment, but when my students have done it, they've built on the foundation you now have.

## Erlang for high-performance computing (really)

- Use existing C++ code for  $\mu\text{-}\varphi$ .
  - It has been carefully optimized – it's fast.
  - It has been widely used over the past 25 years – it's robust.
- Use Erlang to make it parallel!
  - Erlang handles the communication between processes.
  - The code is simple: it works and it's flexible.
  - Erlang can call the C++ functions:
    - The compile intensive part is done in C++.
    - The Erlang code is not a serious bottleneck.

Mark Greenstreet Industrial Strength: Parallel Model Checking CS 418 – Mar. 31, 2017 8 / 17

## Credits

Preventing "overwhelm and crash"

- Drain the inbox into another buffer whenever possible.
- Maintain a credit system
  - When a process  $X$  sends a message to process  $Y$ ,  $X$  decrements its credit-count for  $Y$ .
  - If the credit-count is 0,  $X$  waits to send its message.
  - When  $X$  moves a message from  $Y$  out of its inbox, it sends a credit back to  $X$ .
  - Of course, these messages are piggy-backed on the new-state messages.

Mark Greenstreet Industrial Strength: Parallel Model Checking CS 418 – Mar. 31, 2017 12 / 17

## Summary

PRReach shows how the ideas from this class can be used to build real-world, high-performance, large-scale, parallel systems.

- Lessons learned:
  - Erlang is a great environment for building large-scale, parallel/distributed code.
  - Use the C/C++ call interface to use native C/C++ code for the compute intensive parts of the code.
    - Erlang provides three such interfaces!
  - Watch out for overwhelm and crash
    - If you're going to send a lot of messages, you need some kind of flow control mechanism.

## Memory

- The worklist is the dominant use of memory (in practice)
  - Why?
  - The worklist needs complete state descriptions.
  - The known-state set can use much smaller hashes.
- Solution: store the worklist on disk.
  - Disks are slow – is this crazy?
  - It works just fine because we can access the worklist in any order.
  - Keep a large piece of the worklist in main memory.
  - If the in-memory work list grows too large, then copy a large chunk to disk.
  - The disk reads and writes can be performed asynchronously.
  - See Using magnetic disks ... in the mur... model checker, U. Stern and D.L. Dill.
- Storing the known-state set on disk is much less practical because it's a random-access structure.

Mark Greenstreet Industrial Strength: Parallel Model Checking CS 418 – Mar. 31, 2017 9 / 17

## Load Balancing

- Not all processes have the same amount of work, and they don't all run at the same speed.
- This can lead to **idle processors**.
- Solution:
  - Processes include the length of their worklist in their messages to other workers.
  - If a worker has a short worklist, it asks for half of the worklist of the worker with the longest worklists.
- This is very a coarse-grained approach
  - PRReach makes no effort to keep worklist lengths equal.
  - The coarse-grained approach requires very few messages: **avoid  $\lambda$** .
  - The performance is very good.

## Preview

April 3: Distributed Termination Detection

April 5: Party: 50<sup>th</sup> Anniversary of Amdahl's Law

## Parallel Algorithms 1

### Course Summary

Mark Greenstreet

CpSc 418 – Apr. 5, 2017

## Parallel Algorithms 1

### Course Summary

- map, reduce, and scan: simple patterns
  - Easy to parallelize.
  - Learn to recognize when a problem can be solved by these simple methods.
- sorting networks
  - oblivious computation: the control flow doesn't depend on data values
  - oblivious algorithms are good candidates for parallelism because we can determine the control flow in advance.
  - This lets us identify the data dependencies, and find a parallel solution.
  - The 0-1 principle
  - Bitonic sorting: it's merge sort with an oblivious merge.

Mark Greenstreet Course Summary CS 418 – Apr. 5, 2017 1 / 8

## Parallel Performance

### Parallel Performance

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017

Mark Greenstreet Course Summary CS 418 – Apr. 5, 2017 2 / 8

## Parallel Programming Paradigms

### Parallel Programming Paradigms

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 6 / 8

## Parallel Performance

### Parallel Performance

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 7 / 8

## Parallel Programming Paradigms

### Parallel Programming Paradigms

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 8 / 8

## Parallel Performance

### Parallel Performance

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 9 / 8

## Parallel Programming Paradigms

### Parallel Programming Paradigms

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 10 / 8

## Parallel Performance

### Parallel Performance

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 11 / 8

## Parallel Programming Paradigms

### Parallel Programming Paradigms

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 12 / 8

## Parallel Performance

### Parallel Performance

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 13 / 8

## Parallel Programming Paradigms

### Parallel Programming Paradigms

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 14 / 8

## Parallel Programming Paradigms

### Parallel Programming Paradigms

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 15 / 8

## Parallel Programming Paradigms

### Parallel Programming Paradigms

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 16 / 8

## Parallel Programming Paradigms

### Parallel Programming Paradigms

Mark Greenstreet

Course Summary

CS 418 – Apr. 5, 2017 17 / 8