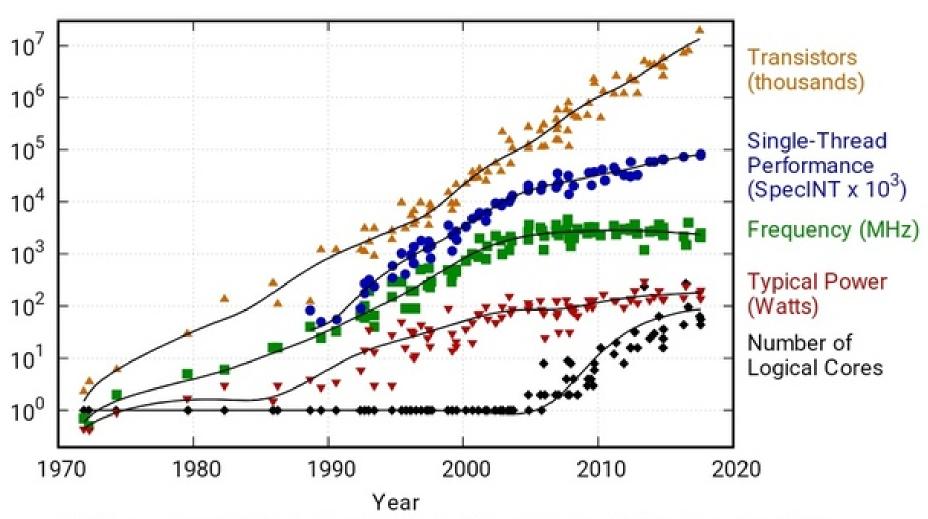
Threads by Compiler

15-411/15-611 Compiler Design

Seth Copen Goldstein

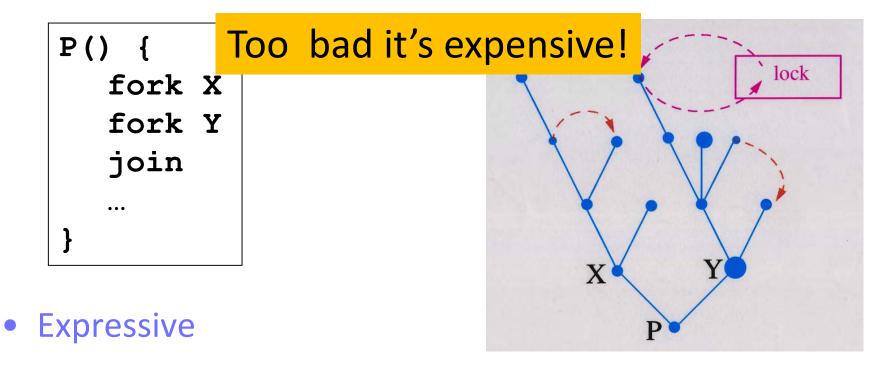
December 3, 2019

The Facts of Life in 2019



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2017 by K. Rupp

The Parallel Call



- Dynamic
- Synchronization dependencies
 - ⇒ must support suspension
- Resource utilization
 - ⇒ must support suspension

There is Hope

Power of call is often not used

Example:

All processors are

Synchronization de

Lock is available

Data request is local

Why is zero-overhead

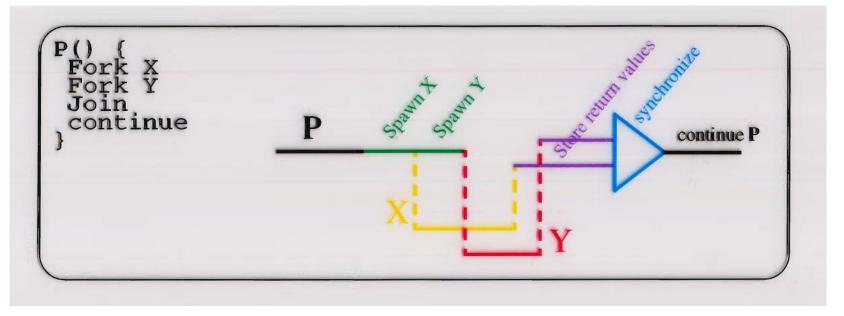
Fork important?

- Goal: zero overhead threading
- More realistically: Pay for what you need.

Parallel Call Inherently Expensive

Parallel Sequential

Parent and child run concurrently

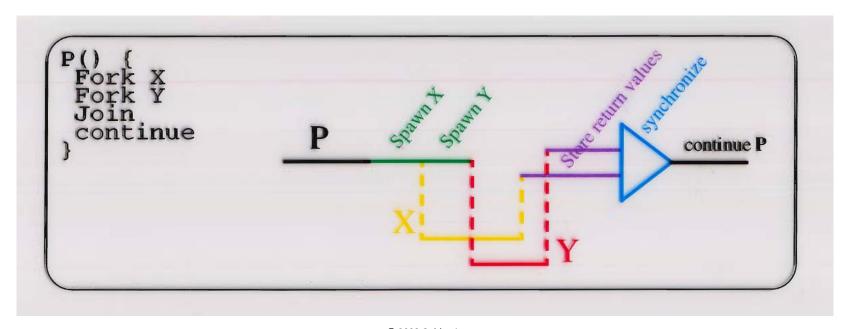


Parallel Call Inherently Expensive

Parallel Sequential

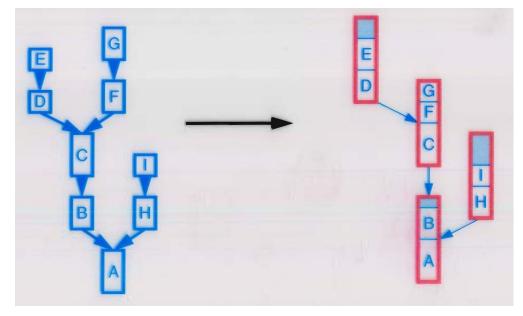
- Parent and child run concurrently
 ⇒ Storage allocated on the heap
- Data and control transferred separately
 - values passed through memory
 - explicit synchronization always needed.

- Parent suspends
 - \Rightarrow use stack
- Transferred together
 - Values in registers
 - No synchronization!

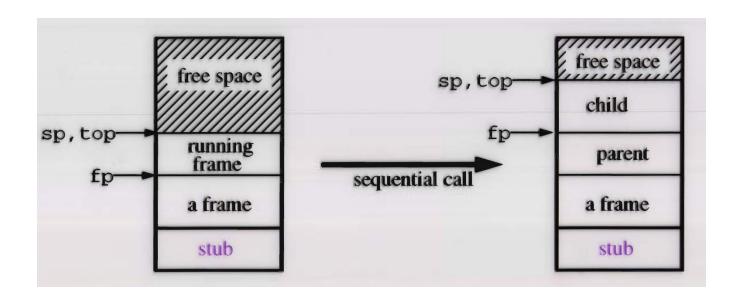


Storage Model

- Cactus Stack
- Invariants:
 - Child returns ⇒ All of its children are done
 - Child suspends ⇒ It currently has no work
 - Parent knows address of child on stack

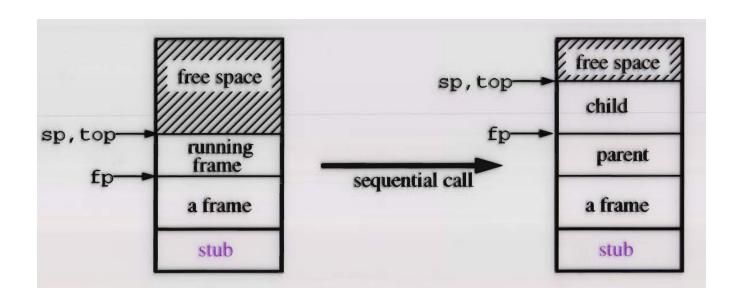


Sequential Call on Stacklet



Overhead?

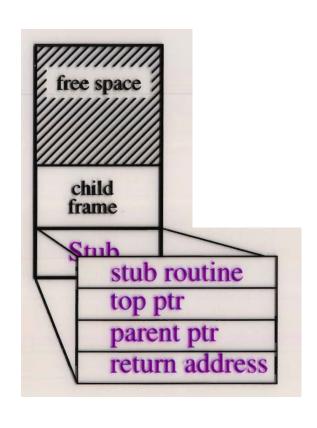
Sequential Call on Stacklet



- Only overhead due to checking for overflow
- Requires changes to the prolog (or, maybe not?)

Stubs

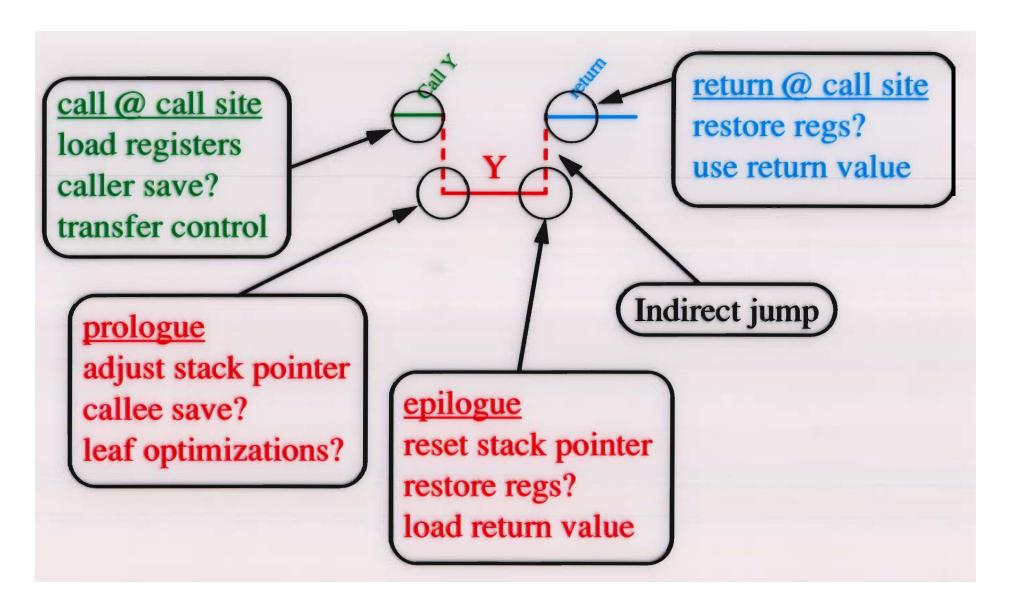
Stub routines handle the special cases



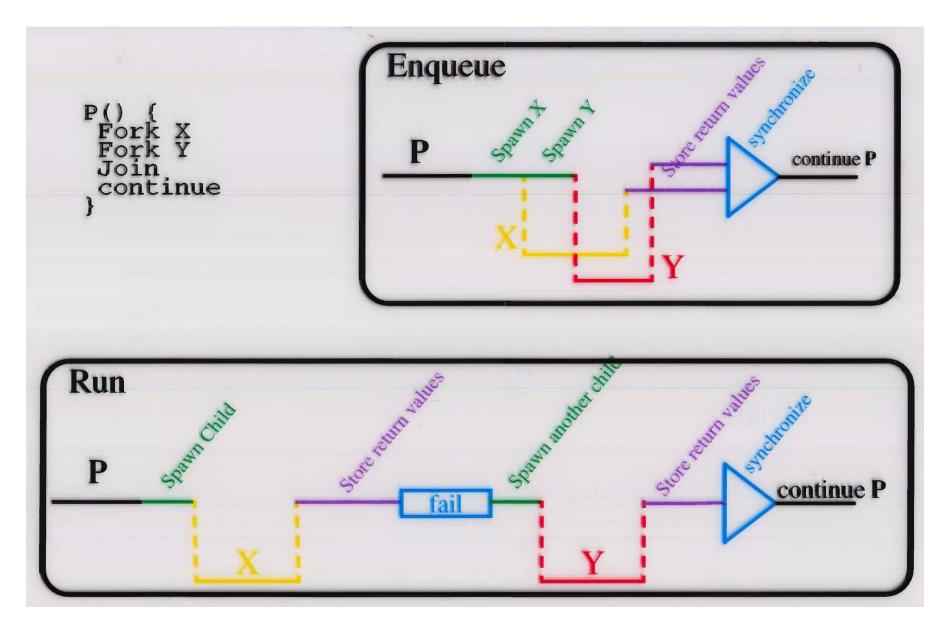
- Stacklet underflow
- Parallel return
- Remote return

Thus, no change in function epilog.

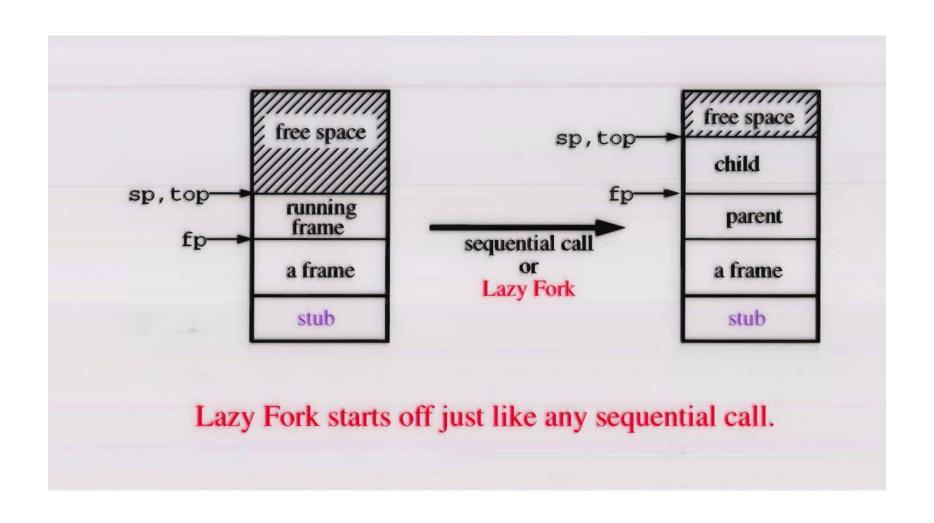
Compiler Focus



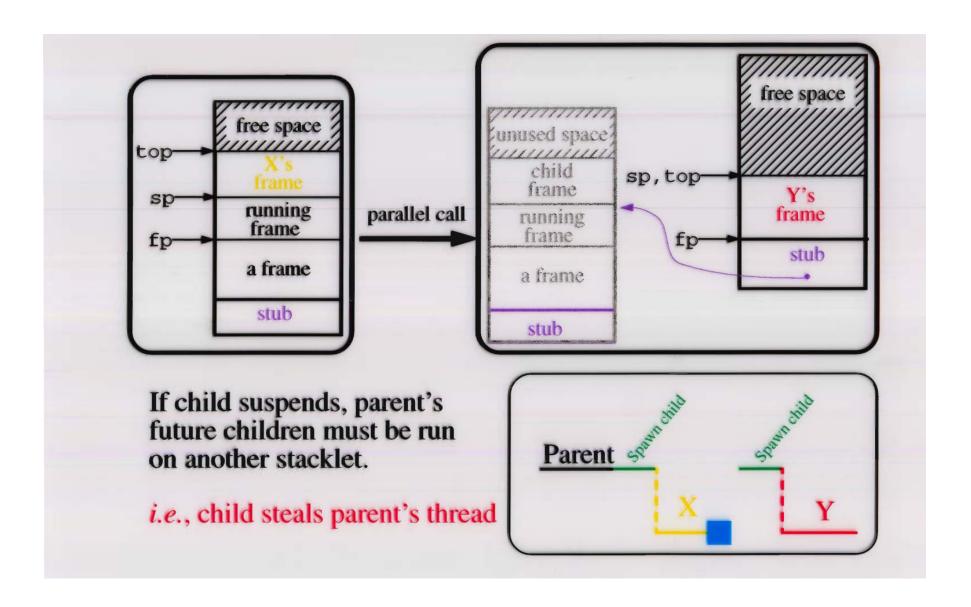
Behavior of parallel call on 1 proc



Lazy Fork



Parallel Allocation

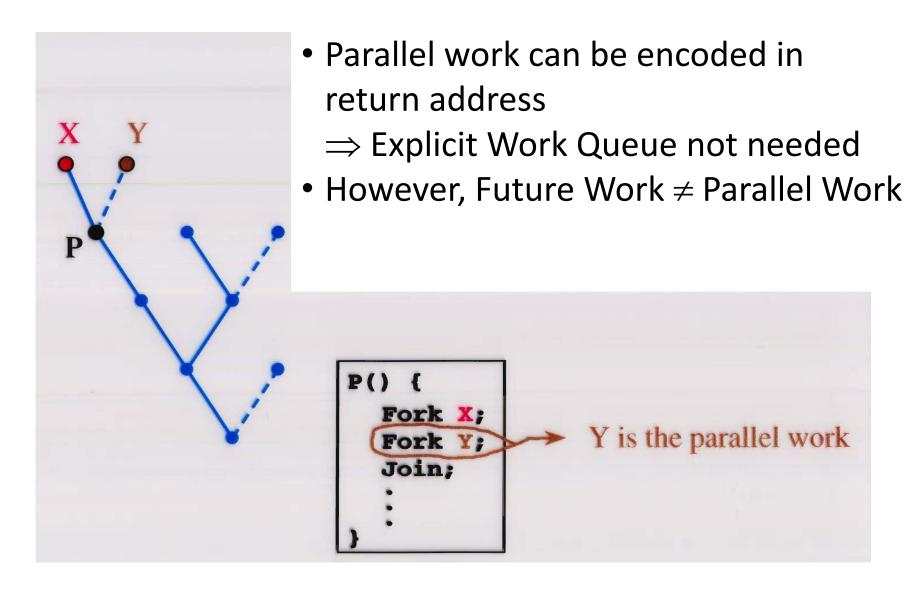


After the Lazy Fork

What if parallelism is required?

- How is child converted into a thread?
- How is remaining work represented?
- How is it located?
- How is it distributed?

Representing Potential Parallel Work



At this point in the execution, Y is ready to execute.

P() {
 Fork X(args);
 Fork Y(args);
 Join;
 :
}

Future Work = Parallel Work

- return address serves 3 purposes:
 - saves results

sequential return

restores parent state

parallelism requires

continues execution at Fork Y

At this point in the execution, there is no parallel work.

```
P() {
   Fork X(args);
   Fork Y(args);
   Join;
   :
}
Parent Fork

Y
```

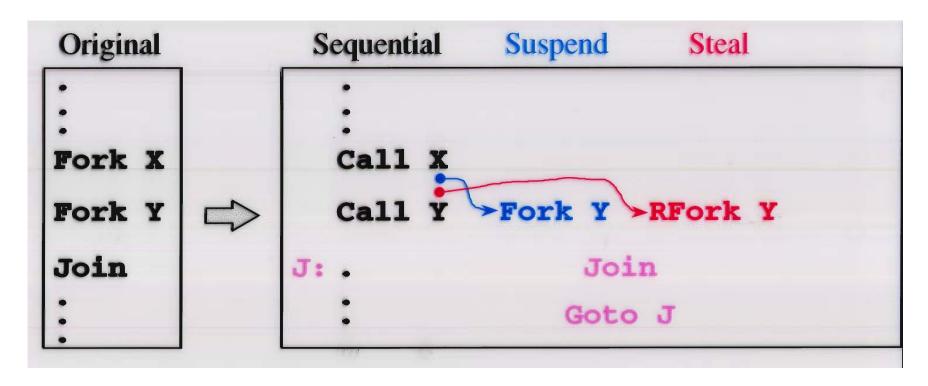
Future Work ≠ Parallel Work

- return address serves 3 purposes:
 - saves results
 - restores parent state
 - continues execution at Join

Thread Seeds

- Use return address to encode
 - Normal return
 - Suspension
 - Work stealing routine
- At time of call, parent plants a seed. Seed is dormant until:
 - Child returns \Rightarrow seed inlined into parent
 - − Child suspends ⇒ seed is activated
 - Remote core grabs work
 - \Rightarrow seed is stolen

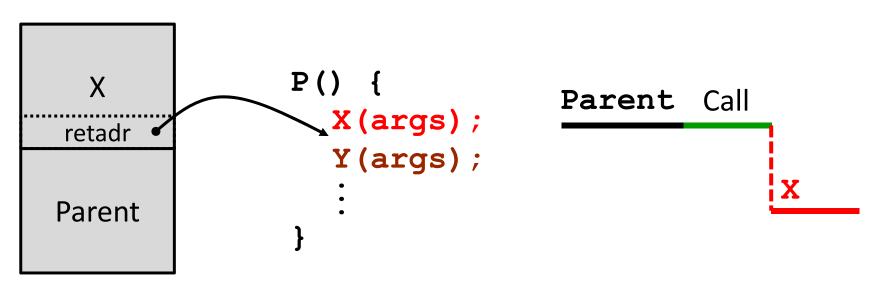
Code Gen Strategy



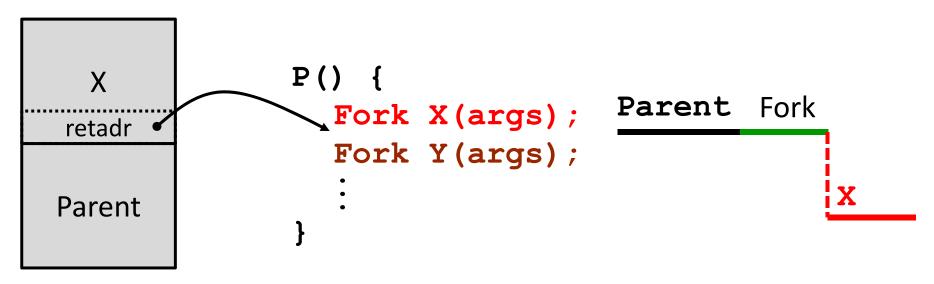
Create Three code streams:

- + Eliminates unnecessary synchronization
- + Eliminates unnecessary bookkeeping
- + Optimizes for common case
- Increases code size

- Must encode state of parent:
 - Never had parallel work.



- Must encode state of parent:
 - Never had parallel work.
 - Has potentially parallel work.
 (currently has sequentially invoked child)



- Must encode state of parent:
 - Never had parallel work.
 - Has potentially parallel work.
 (currently has sequentially invoked child)
 - Has parallel work.
 (currently has converted sequentially invoked child)

```
//save result
if (--synch == 0) enq(P);
return;

P() {
   Fork X(args);
   Fork Y(args);
    :
}
Parent

Y

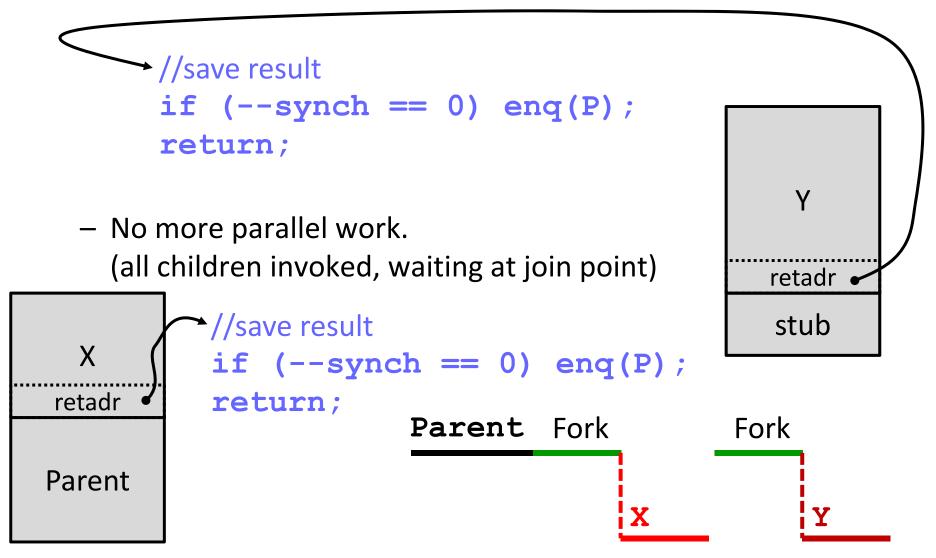
Parent

Y

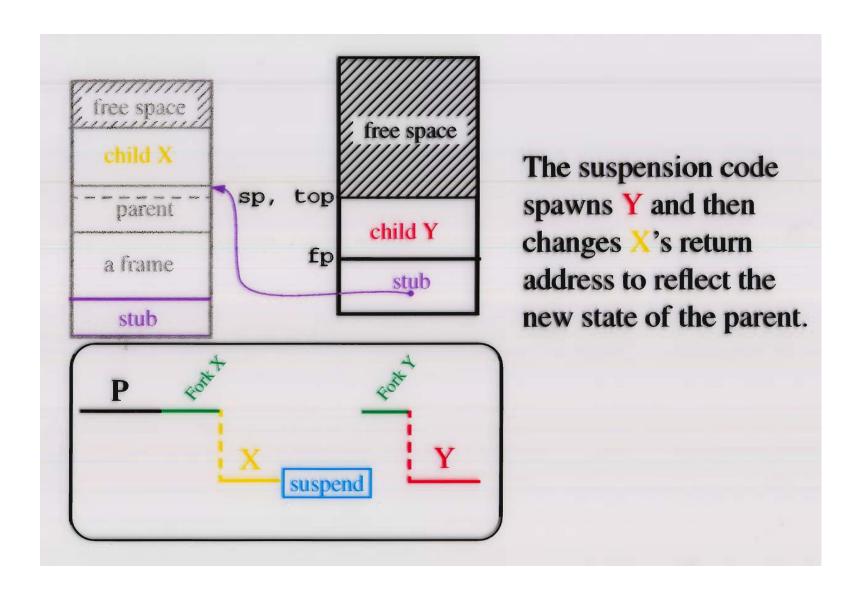
Parent

Fork
Y(args);
X
```

Must encode state of parent:



Parent Controlled Return Continuations



Dealing with state

In sequential context want a, b, c, d, ...
 mapped to registers.

```
P() {
    c = Fork X(a);
    d = Fork Y(b);
    Join;
    :
}
```

Dealing with state

In sequential context want a, b, c, d, ...
 mapped to registers.

```
P() {
    c = Fork X(a);
    d = Fork Y(b);
    Join;
    :
}
```

```
P() {
  c = call X(a);
  d = call Y(b);
  // Join;
  // converted X ret
  c ← return value
  sync code
  // spawn Y
  change X's retadr
  Fork Y(b);
  d ← return value
  synch code
```

Dealing with state

In sequential context want a, b, c, d, ...
 mapped to registers.

```
P() {
    c = Fork X(a);
    d = Fork Y(b);
    Join;
    :
}
```

However, if converted,
 Fork Y needs access to b and return from X needs to access c.

```
P() {
  c = call X(a);
  d = call Y(b);
  // Join;
   // converted X ret
  c \leftarrow return value
  sync code
  // spawn Y
  change X's retadr
  Fork Y(b);
  d \leftarrow return value
  synch code
```

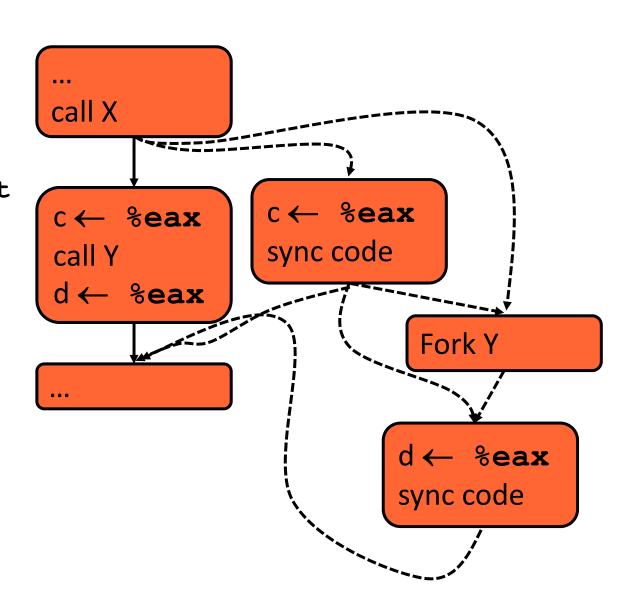
Original CFG

```
P() {
  c = call X(a);
  d = call Y(b);
  // Join;
   // converted X ret
  c ← return value
  sync code
  // spawn Y
  change X's retadr
  Fork Y(b);
  d ← return value
  synch code
```

```
call X
c ← %eax
call Y
d ← %eax
//Join
....
```

Modified CFG

```
P() {
  c = call X(a);
  d = call Y(b);
  // Join;
   // converted X ret
  c ← return value
  sync code
  // spawn Y
  change X's retadr
  Fork Y(b);
  d ← return value
  synch code
```



Modified CFG

```
P() {
                                          In reg
  c = call X(a);
  d = call Y(b);
                                                  In frame
                         call X
   // Join;
   // converted X ret
                                       c ← %eax
                         c ← %eax
  c ← return value
                                       sync code
  sync code
  // spawn Y
                         d ← %eax
  change X's retadr
                                                  Fork Y
  Fork Y(b);
  d ← return value
  synch code
                                                d ← %eax
                                                sync code
                           Reload
                           registers
```

Implementing the Seed

- Every return address has an entry in a per procedure table mapping the retadr to:
 - conversion routine
 - local fork routine
 - remote fork routine
 - Information about stack size, registers, etc.
- Runtime looks at top of stack
 - locates most recent return address
 - finds entry in table
 - executes appropriate code using appropriate frame pointer

Compiled setjmp, longjmp

```
P() {
  c = call X(a);
  d = call Y(b);
  // Join;
  // converted X ret
  c ← return value
  sync code
  // spawn Y
  change X's retadr
  Fork Y(b);
  d ← return value
  synch code
```

? retadr X retadr Parent

Synchronization

- There is a big data race: the return address!
- Locking on every call and return violates "pay for what you use"
- Possible Solutions:
 - user-level interrupts and message passing
 - fancy tricks with atomic xchg and do explicit indirect jmp for return

Summary

- Key: Only create parallel work when necessary.
- Exposing threads to the compiler is essential for high-performance fine-grained multi-threading.
- There are three axis to consider:
 - How to store thread state
 - Representation of (potentially) parallel work
 - How to convert potentially parallel work to actual parallel work
- Exploit compiler
- Exploit indirect jump

Logistics

- Lab6 due next week
- NO EXTENSIONS!
- FCEs
- More detailed feedback

A very big Thank-you!