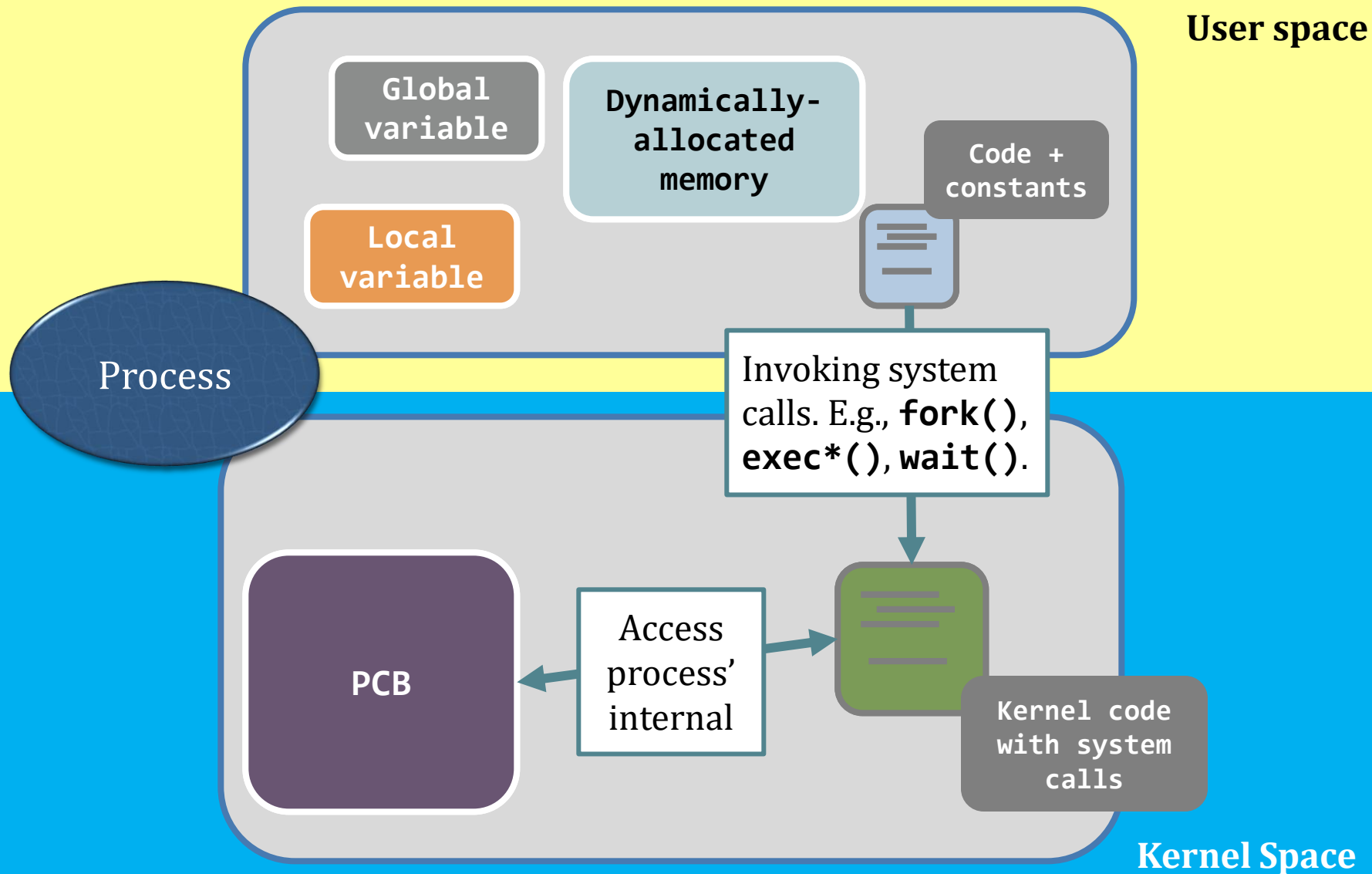# Lecture 4: Process II
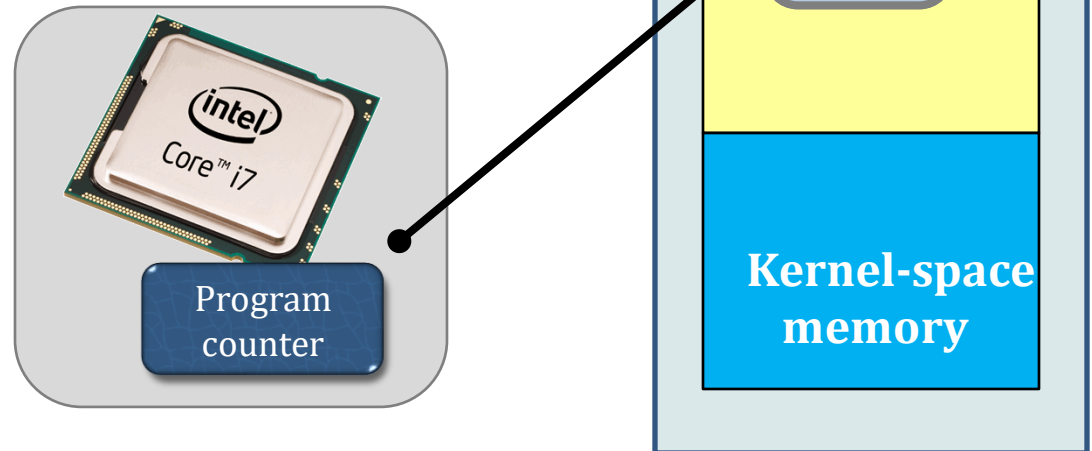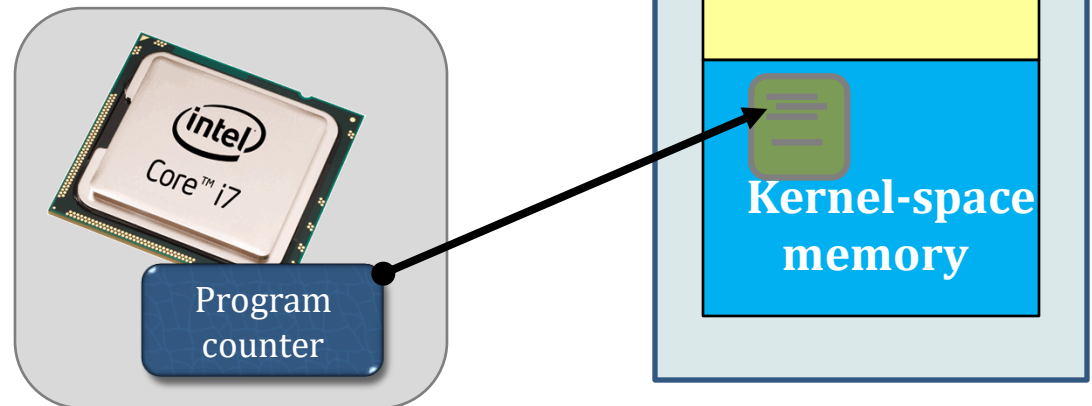
Bo Tang @ 2020, Spring

# The story so far…

# When invoking a system call (memory view)

◈ When running a program code of a user process.

◈ As the code is in user-space memory, so the program counter is pointing to that region.

**Memory**

User-space memory

Kernel-space memory

Program counter

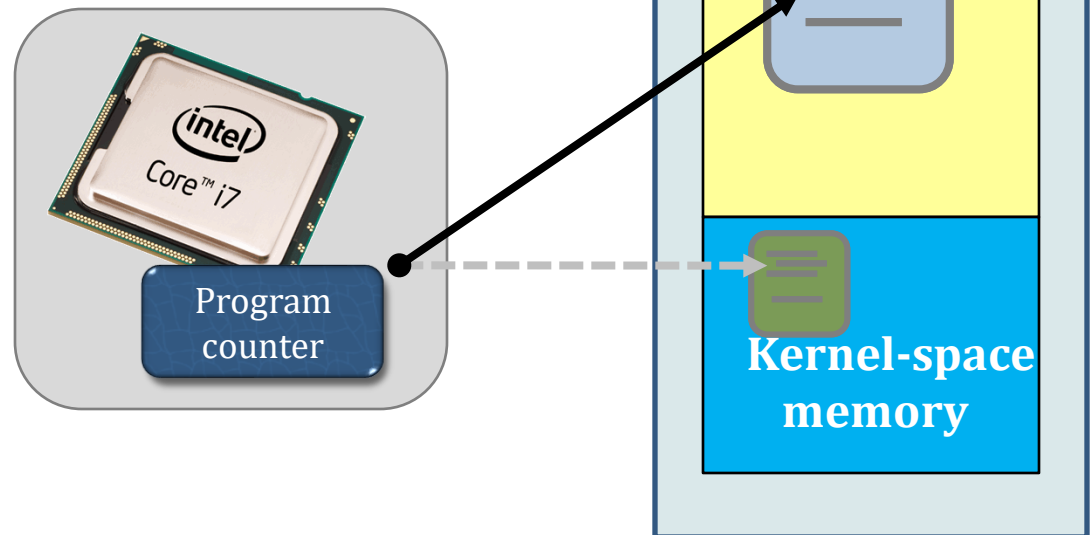# When invoking a system call (memory view)

◈ When the process is calling the system call "**getpid()**".

◈ Then, the CPU switches <u>from the user-space to the kernel-space</u>, and reads the PID of the process from the kernel.

**Memory**

**User-space memory**
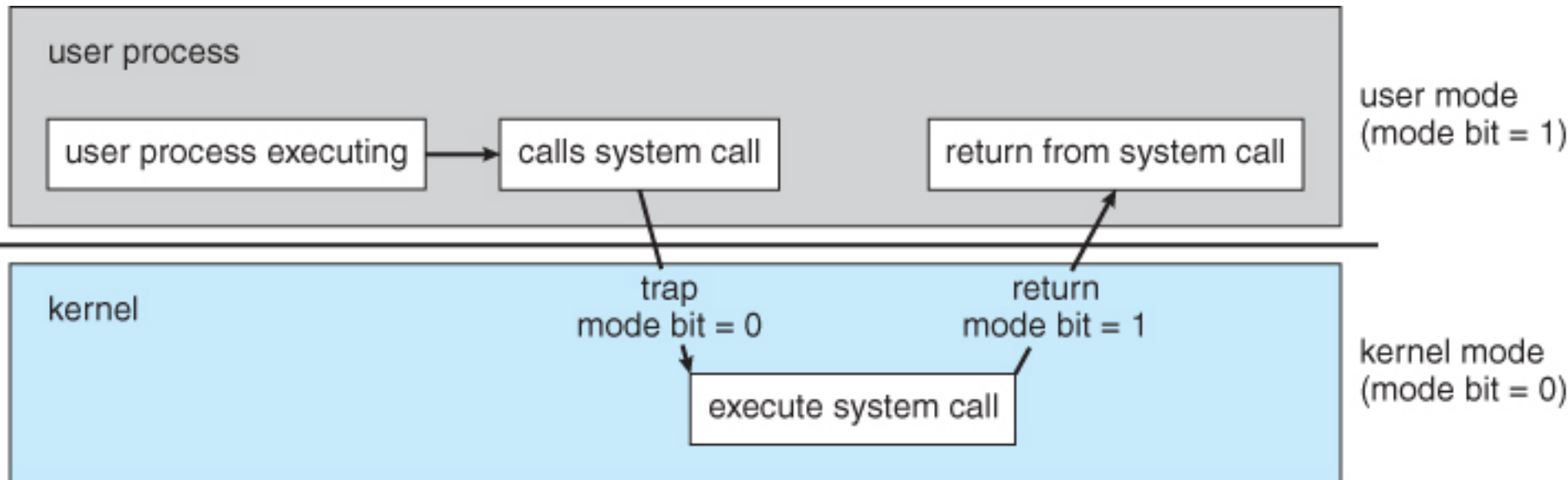
**Kernel-space memory**

Program counter

# When invoking a system call (memory view)

◈ When the CPU has finished executing the "**`getpid()`**" system call

  ◈ it <u>switches back to the user-space memory</u>, and continues running that program code.

**Memory**

User-space memory

Kernel-space memory

Program counter

# When invoking a system call (CPU view)

# Process real time cost (wall-clock time)

calling system call.
  e.g., **getpid()**

Some system calls may take a long time.
E.g., accessing a printer.

**User time** –
CPU time spent on codes in user-space memory.

**Sys time** –
CPU time spent on codes in kernel-space memory.

Read information and the system call returns.

# User time  VS  System time – example 1

◈ Let's tell the difference…with the tool "`time`".

```
$ time ./time_example

real     0m0.001s
user     0m0.000s
sys      0m0.000s
$ _
```

**Real-time** elapsed when "`./time_example`" terminates.

The **user time** of "**./time_example**".

The **sys time** of "**./time_example**".

It's possible:
real > user + sys
real < user + sys

Why?

```c
int main(void) {
    int x = 0;
    for(i = 1; i <= 10000; i++) {
        x = x + i;
//      printf("x = %d\n", x);
    }
    return 0;
}
```

# User time VS System time – example 1

◈ Let's tell the difference…with the tool "`time`".

```
$ time ./time_example

real      0m0.001s
user      0m0.000s
sys       0m0.000s
$ _
```

```
int main(void) {
    int x = 0;
    for(i = 1; i <= 10000; i++) {
        x = x + i;
        // printf("x = %d\n", x);
    }
    return 0;      Commented on purpose.
}
```

```
$ time ./time_example

real 0m2.795s
user 0m0.084s
sys  0m0.124s
$ _
```

See? Accessing hardware costs the process more time.

```
int main(void) {
    int x = 0;
    for(i = 1; i <= 10000; i++) {
        x = x + i;
        printf("x = %d\n", x);
    }
    return 0;      Comment released.
}
```

# User time VS Sys time – example 2

◈ The user time and the sys time together **define the performance of an application**.

◈ When writing a program, you must consider both the user time and the sys time.

◆ E.g., the output of the following two programs are exactly the same. But, their running time is not.

```
#define MAX 1000000

int main(void) {
    int i;
    for(i = 0; i < MAX; i++)
        printf("x\n");
    return 0;
}
```

```
#define MAX 1000000

int main(void) {
    int i;
    for(i = 0; i < MAX / 5 ; i++)
        printf("x\nx\nx\nx\nx\n");
    return 0;
}
```

# User time VS Sys time – example 2

```
#define MAX 1000000

int main(void) {
    int i;
    for(i = 0; i < MAX; i++)
        printf("x\n");
    return 0;
}
```

```
$ time ./time_example_slow

real 0m1.562s
user 0m0.024s
sys  0m0.108s
$ _
```

```
#define MAX 1000000

int main(void) {
    int i;
    for(i = 0; i < MAX / 5 ; i++)
        printf("x\nx\nx\nx\nx\n");
    return 0;
}
```

```
$ time ./time_example_fast

real 0m1.293s
user 0m0.012s
sys  0m0.084s
$ _
```

# User time  VS  Sys time

- Function calls cause overhead
  - Stack pushing (will see later)
- Sys calls may cause even more
  - ➔ Sys call is from another "process" (the kernel)
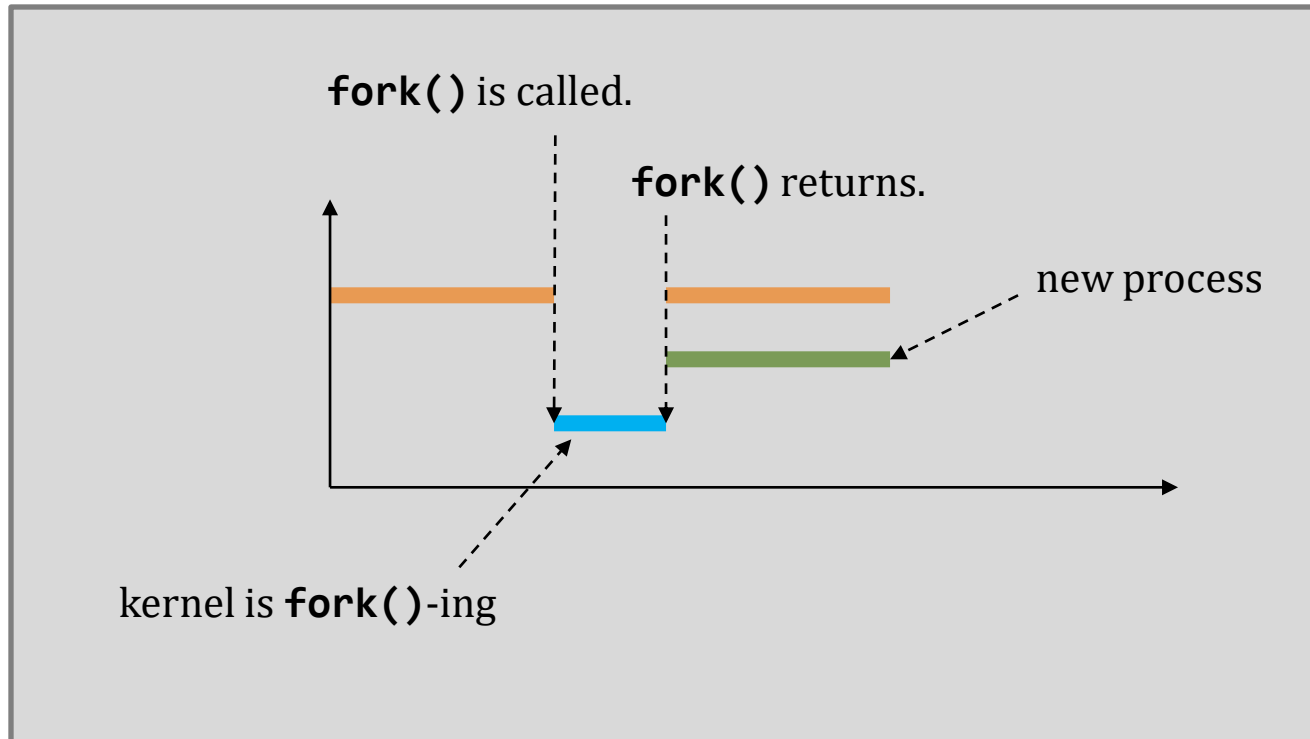  - ➔ Switching to another "process" ➔ context switch (will see later)

https://www.quora.com/Is-an-OS-kernel-itself-a-process

# Working of system calls
##    - fork();

**Process**

# Programmer view of fork()

fork() is called.

fork() returns.

new process

kernel is fork()-ing

# **fork()** inside the kernel

This guy invoked **fork()**.

**OS Kernel**

Inside kernel, processes are arranged as a doubly linked list, called the *task list*.

| Process 1234 | Process 345 | ● ● ● |

| PID = 1234 |
|---|
| Running time |
| Array of opened files |
| ● ● ● |

**copying**

| PID = 1234 |
|---|
| Running time |
| Array of opened files |
| ● ● ● |

# **fork()** in action – kernel-space update

# **fork()** in action – user-space update

# **fork()** in action – finish



**OS Kernel**

Ready to return from **fork()**

Ready to return from **fork()**

Process 1234 → Process 345 → Process 1235 → • •

| PID = 1234 |
| --- |
| Running time |
| Array of opened files |
| List of children |
| Return value = 1235 |

| PID = 1235 |
| --- |
| Running time |
| Array of opened files |
| Pointer to my parent |
| Return value = 0 |

18

# **fork()** in action – array of opened files?

◈ Array of opened files contains:

| Array Index | Description |
|---|---|
| 0 | Standard Input Stream; **FILE *stdin;** |
| 1 | Standard Output Stream; **FILE *stdout;** |
| 2 | Standard Error Stream; **FILE *stderr;** |
| 3 or beyond | Storing the files you opened, e.g., **fopen()**, **open()**, etc. |

◈ That's why a parent process **shares the same terminal output stream** as the child process.

Stream is just a **logical** object for you to read as a **sequence of bytes**
So, how can you random access the middle of a file?  Read Stream → Array → Array[mid-point]

# **fork()** in action – sharing opened files?

◈ What if two processes, **sharing the same opened file**, write to that file together?

Process    Process

write()    write()

Let's see what will happen when the program finishes running!

# Working of system calls

- `fork();`
- `exec*();`

**Process**

# exec*() that you've learnt...

◈ How about the **exec*()** call family?

**exec*()** is called.

Old code → New code

**Process**

The process returns to user-space **but is executing another program**.

What is the kernel is doing?

# exec*() in action

This guy invoked **exec*().**

**OS Kernel**

**Process 1234** → **Process 345** → • • •

The kernel will also
- reset the register values (e.g., program counter)

Cleared!

**Local variable**

**Dynamically-allocated memory**

Cleared!

Reset based on the new code!

**Global variable**

**Code + constants**

Changed to the new program code!

User space

# Working of system calls

```
- fork();
- exec*();
- wait() + exit();
```

Process

Process

# wait() and exit()



wait() is called.

wait() blocks the parent.

wait() returns.

Parent

Child

**Parent Process**

**Child Process**

Child is terminated through the exit() system call.

How do the two processes communicate?

# exit() (kernel-view)

OS Kernel

This guy invoked **exit()**.

Parent

Child

Process 1234

Process 1235

• • •

PID = 1235

Running time

Array of opened files

• • •

The kernel frees all the allocated memory.

The list of opened files are all closed. (so it's okay to skip **fclose()**; though not recommended)

# exit() (kernel-view)

OS Kernel

This guy invoked **exit()**.

Parent

Child

**Process 1234**

**Process 1235**

• • •

Then, the kernel **frees everything on the user-space memory** about the concerned process, including program code and allocated memory.

Local variable

Dynamically-allocated memory

Global variable

Code + constants

User space

# exit() (kernel-view)

**OS Kernel**

This guy invoked **exit()**.

**Parent**

**Child**

Process 1234

Process 1235

• • •

Process ID stills in the kernel's process table
- Why?
  [Wiki] *This entry is still needed to allow the process that started the (now zombie) process to read its exit status.*
The status of the child is now called **zombie ("terminated")**.

PID = 1235
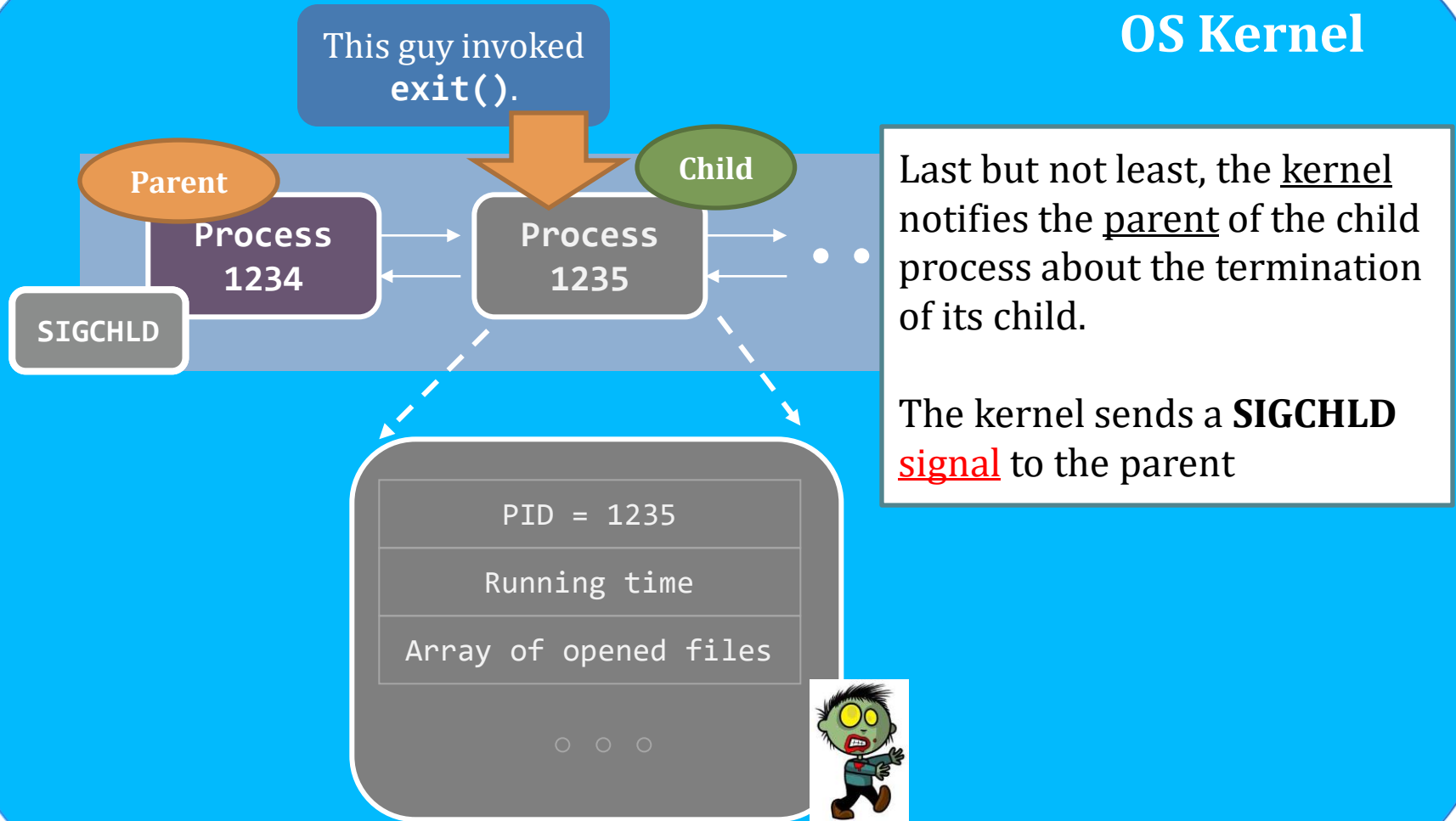
Running time

Array of opened files

○  ○  ○

# exit() (kernel-view)



**OS Kernel**

This guy invoked `exit()`.

Parent

Child

**Process 1234**

**Process 1235**

**SIGCHLD**

Last but not least, the <u>kernel</u> notifies the <u>parent</u> of the child process about the termination of its child.

The kernel sends a **SIGCHLD** <u>signal</u> to the parent

PID = 1235

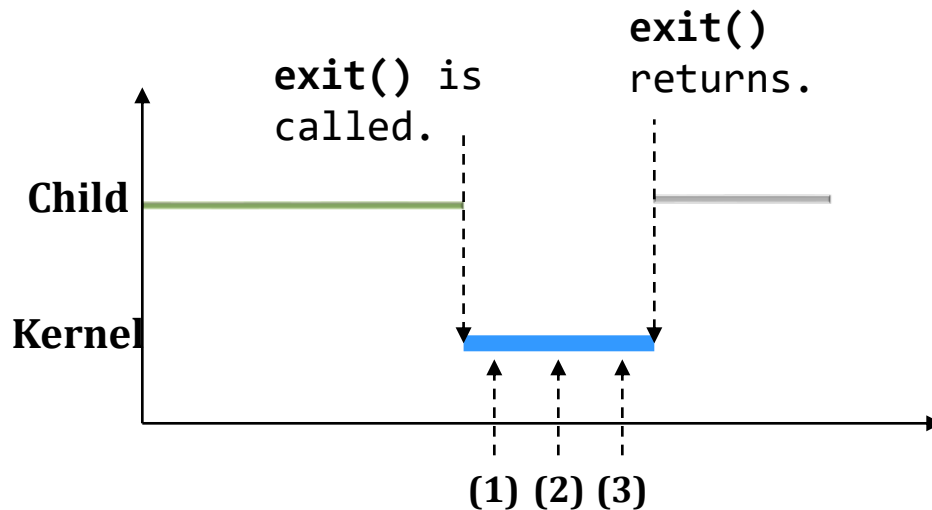Running time

Array of opened files

○ ○ ○

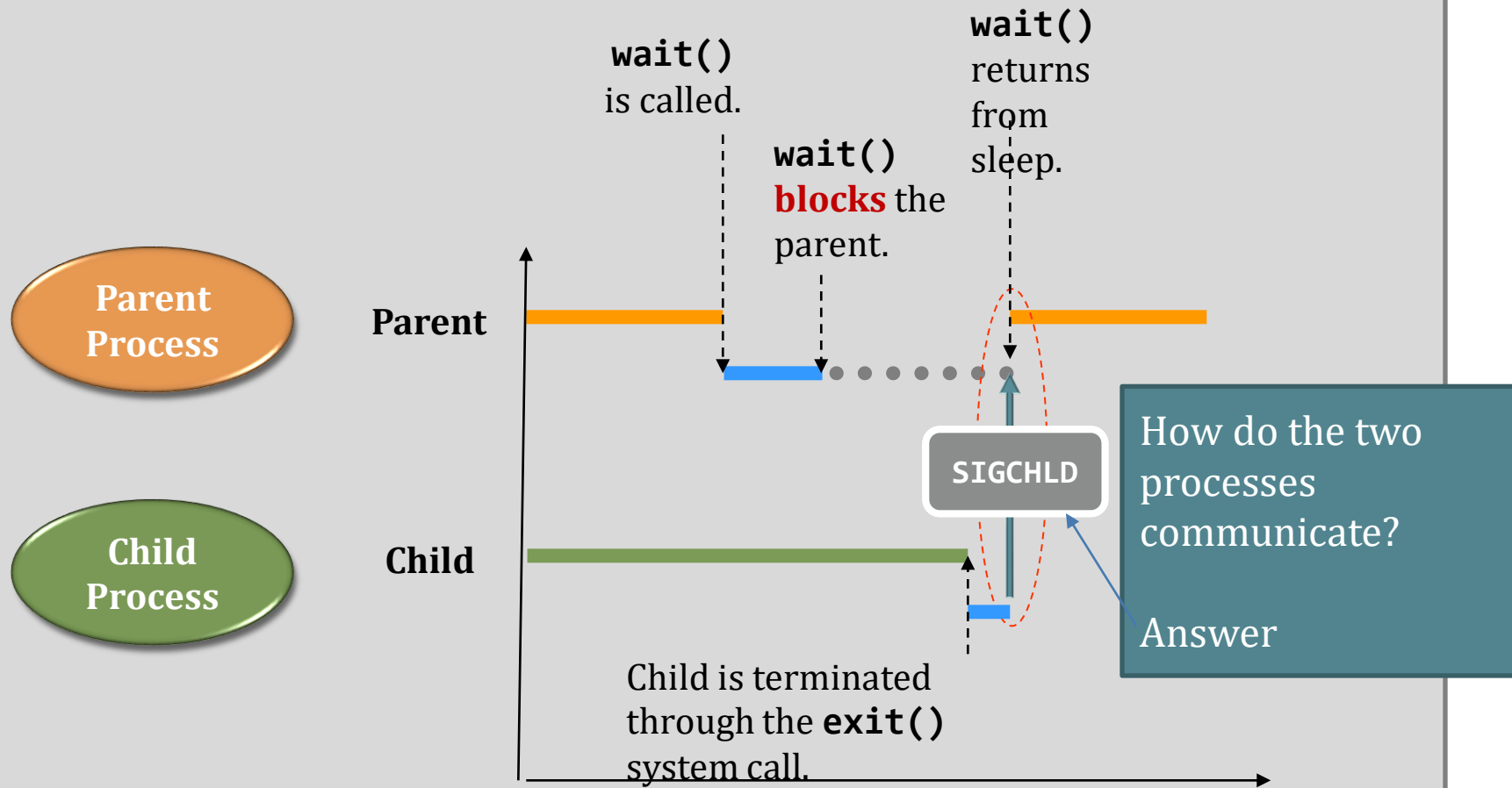# Summary -- what the kernel does for **exit()**

Step (1) Clean up most of the allocated kernel-space memory (e.g., process's running time info).

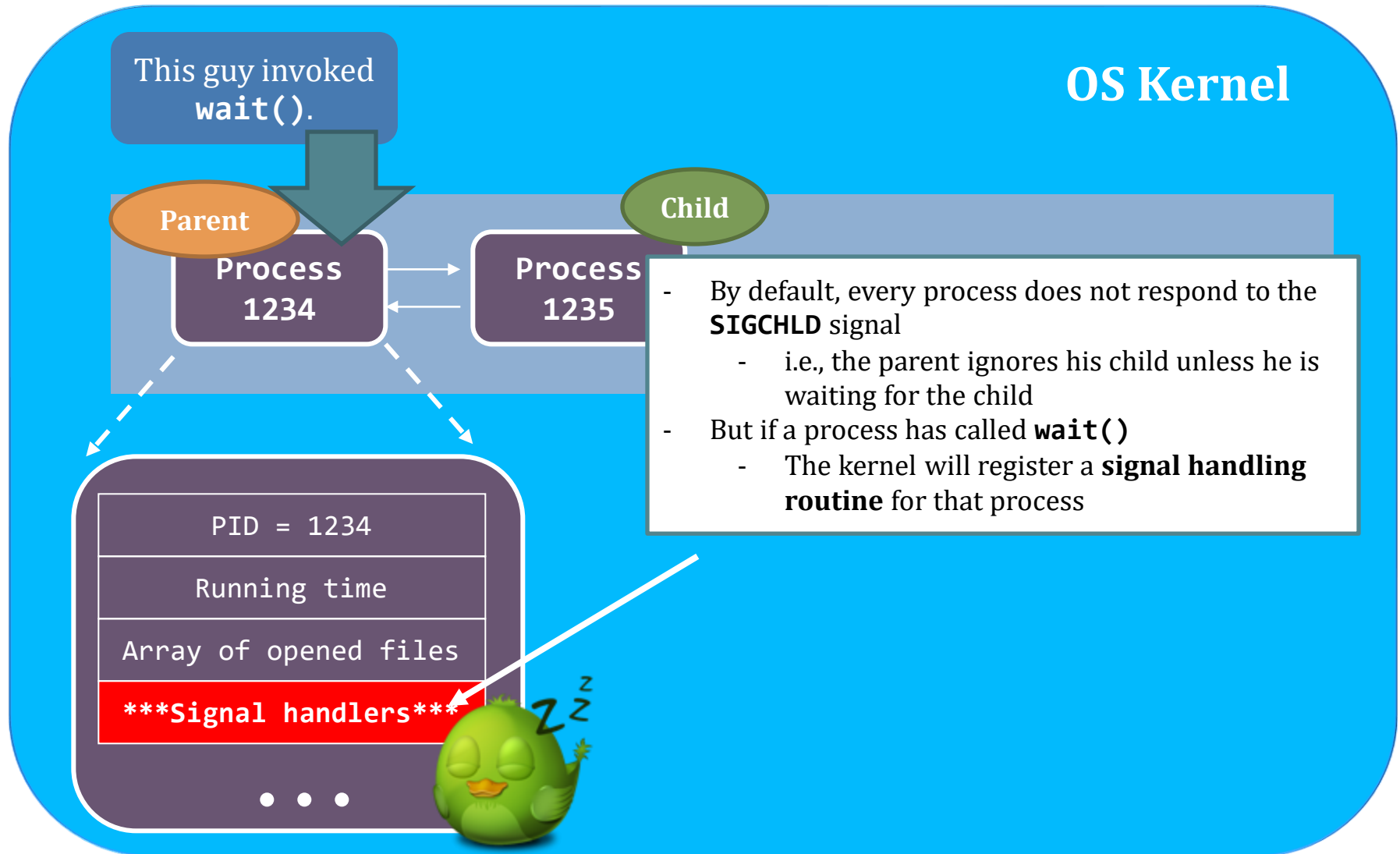Step (2) Clean up the exit process's user-space memory.

Step (3) Notify the parent with SIGCHLD.

**exit()** is
called.

**exit()**
returns.

**Child**

**Kernel**

**(1) (2) (3)**

# wait() and exit()

# **wait()** kernel view's – registering signal handling routine

## OS Kernel

This guy invoked **wait()**.

**Parent**

**Child**

**Process 1234**

**Process 1235**

- By default, every process does not respond to the **SIGCHLD** signal
  - i.e., the parent ignores his child unless he is waiting for the child
- But if a process has called **wait()**
  - The kernel will register a **signal handling routine** for that process

| PID = 1234 |
| Running time |
| Array of opened files |
| ***Signal handlers*** |
| ● ● ● |

# **wait()** kernel's view

This guy invoked **wait()**.

**Parent**

**Child**

**Process 1234**

**Process 1235**

When **SIGCHLD** comes, the corresponding **signal handling routine is invoked**!

Note:
- the parent is still inside the **wait()** system call

PID = 1234

Running time

Array of opened files

Signal handlers

• • •

SIGCHLD from 1235

# wait() kernel's view

# `wait()` kernel's view

OS Kernel

Ready to return from `wait()`.

Parent

Process 1234

The kernel
- deregisters the **signal handling routine** for the parent
- returns the PID of the terminated child as the return value of `wait()`

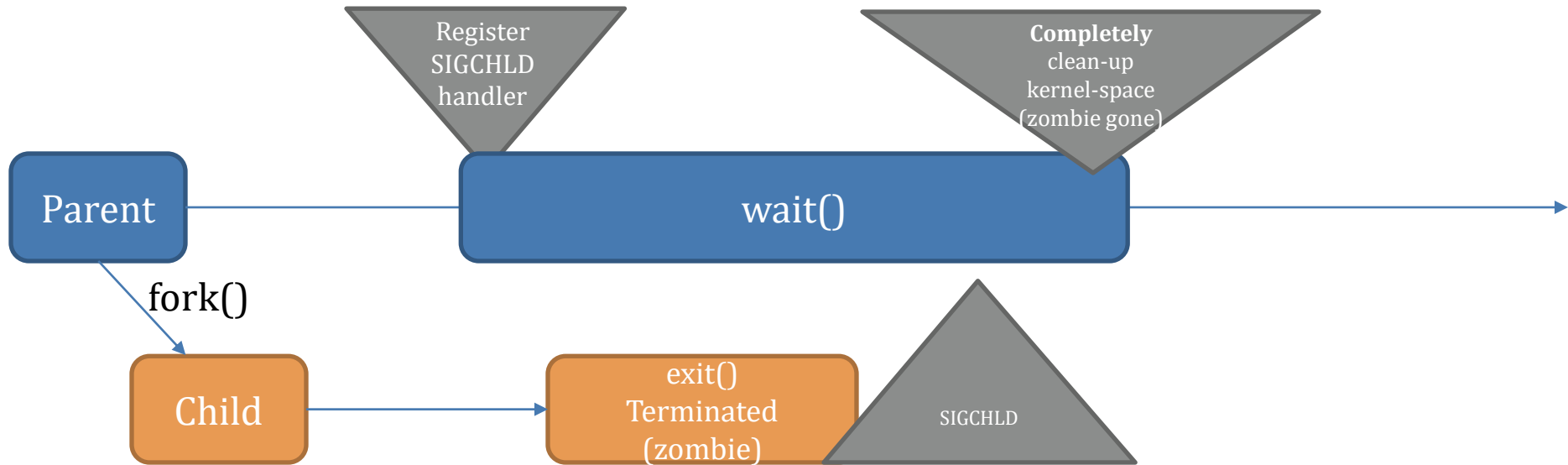The parent is ignoring **SIGCHLD** again.

| PID = 1234 |
| Running time |
| Array of opened files |
| ~~Signal handlers~~ |
| Return value = 1235 |

# Overall – normal case

# Normal Case



**wait()** is called.

**wait()** **blocks** the parent.

**wait()** returns.

Parent

SIGCHLD

Child

Child is terminated through the **exit()** system call.

So, the child will be given a clean death because of the **wait()** system call.

# Parent's wait() after Child's exit()

**Parent Process**

**Child Process**

Parent

Child

**wait()** is called.

… busy...while-loop ...

SIGCHLD

Child is terminated through the **exit()** system call.

# Parent's Wait() after Child's exit()



**OS Kernel**

This guy invoked **wait()**.

**Parent**

Process 1234

**Child**

Process 1235

PID = 1234

Running time

Array of opened files

Signal handlers

SIGCHLD from 1235

On Parent's wait(), kernel:
- sets the signal handling routine
- once set, found SIGCHLD is already there
- takes action immediately

# Parent's Wait() after Child's exit()

This case is okay but the zombie wanders around until **wait()** is called

**wait()** is called.

**wait()** returns.

**Parent Process**

**Parent**

SIGCHLD

**Child Process**

**Child**

Child is terminated through the **exit()** system call.

# **wait()** and **exit()** – short summary

◈ **exit()** system call turns a process into a zombie when...

  ◈ The process calls **exit()**.

  ◈ The process returns from **main()**.

  ◈ The process terminates abnormally.

    ◆ The kernel knows that the process is terminated abnormally. Hence, the kernel invokes **exit()** for it.
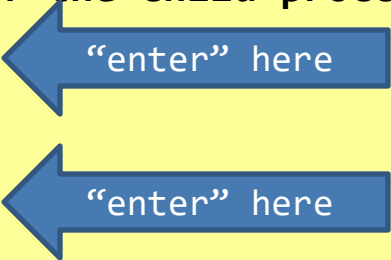
# **wait()** and **exit()** – short summary

◈ **wait()** & **waitpid()** are to reap zombie child processes.

    ◈ It is a must that you should never leave any zombies in the system.

    ◈ **wait()** & **waitpid()** pause the caller until

        ◆ A child terminates/stops, OR

        ◆ The caller receives a signal (i.e., the signal interrupted the `wait()`)

◈ Linux will label zombie processes as "**<defunct>**".

    ◈ To look for them:

```
$ ps aux | grep defunct
….....     3150  ...  [ls] <defunct>
$ _
```

PID of the process

# **wait()** and **exit()** – short summary

```
1  int main(void)
2  {
3      int pid;
4      if( (pid = fork()) !=0 ) {
5          printf("Look at the status of the child process %d\n", pid);
6          while( getchar() != '\n' );        "enter" here
7          wait(NULL);
8          printf("Look again!\n");
9          while( getchar() != '\n' );        "enter" here
10     }
11     return 0;
12 }
```

This program requires you to type "enter" twice before the process terminates.

You are expected to see **the status of the child process changes (ps aux [PID])** between the 1st and the 2nd "enter".
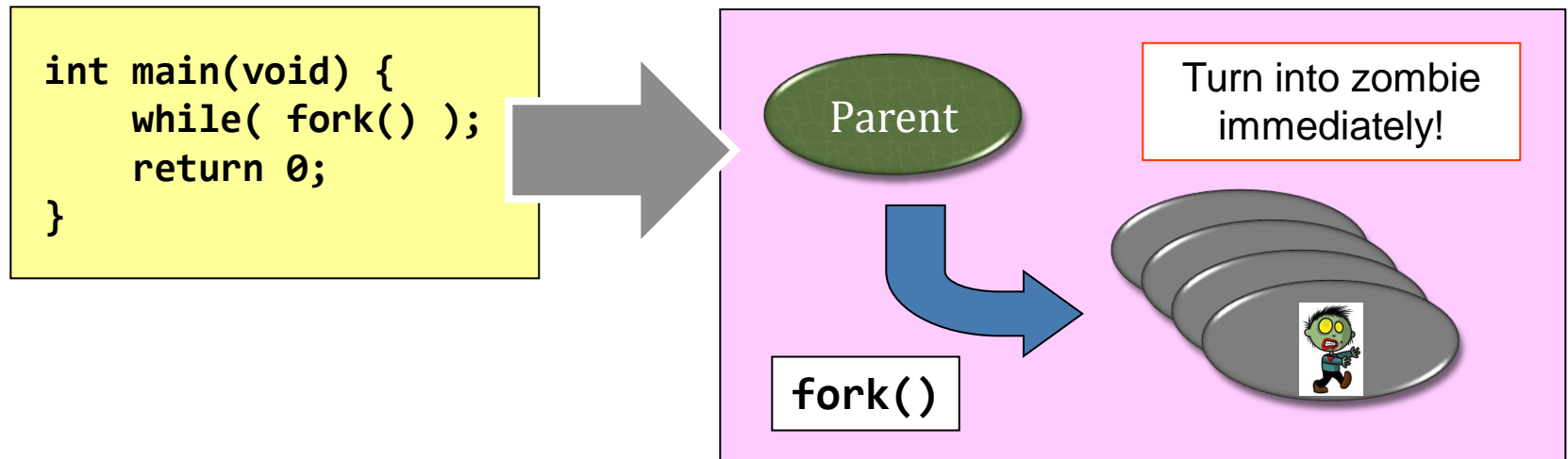
# Working of system calls
```
    - fork();
    - exec*();
    - wait() + exit();
    - importance/fun in knowing
      the above things?
```

# Calling **wait()** is important.

◈ It is not only about process execution / suspension…

◈ It is about **system resource management**.

  ◈ A zombie takes up a PID;

  ◈ The total number of PIDs are limited;

  ◆ Read the limit: "`cat /proc/sys/kernel/pid_max`"

  ◆ It is 32,768.

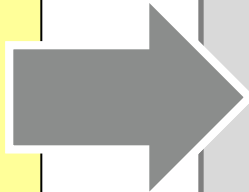  ◈ What will happen if we don't clean up the zombies?

# The fork bomb

◈ Deliberately missing wait()

◈ Do not try this on department's machines…

```
int main(void) {
    while( fork() );
    return 0;
}
```

Parent

Turn into zombie immediately!

fork()

**An infinite, zombie factory!**

# When **wait()** is absent…

```
int main(void) {
    while( fork() );
    return 0;
}
```

```
$ ./interesting

_
                          Terminal A
```

```
$ ls
No process left.
$ poweroff
No process left.
$ =__=
No process left.
$ _
                          Terminal B
```
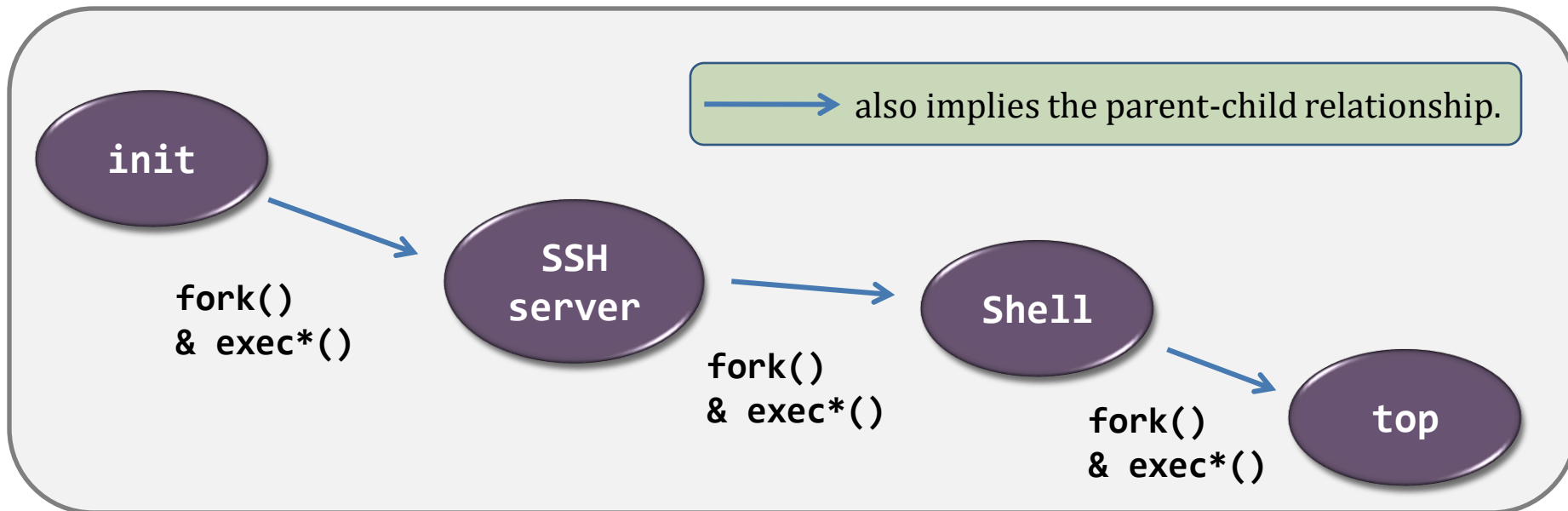
# The first process

- We now focus on the process-related events.
  - The kernel, while it is booting up, creates the first process – **init**.

- The "**init**" process:
  - has **PID = 1**, and
  - is running the program code "**/sbin/init**".

- Its first task is to **create more processes**…
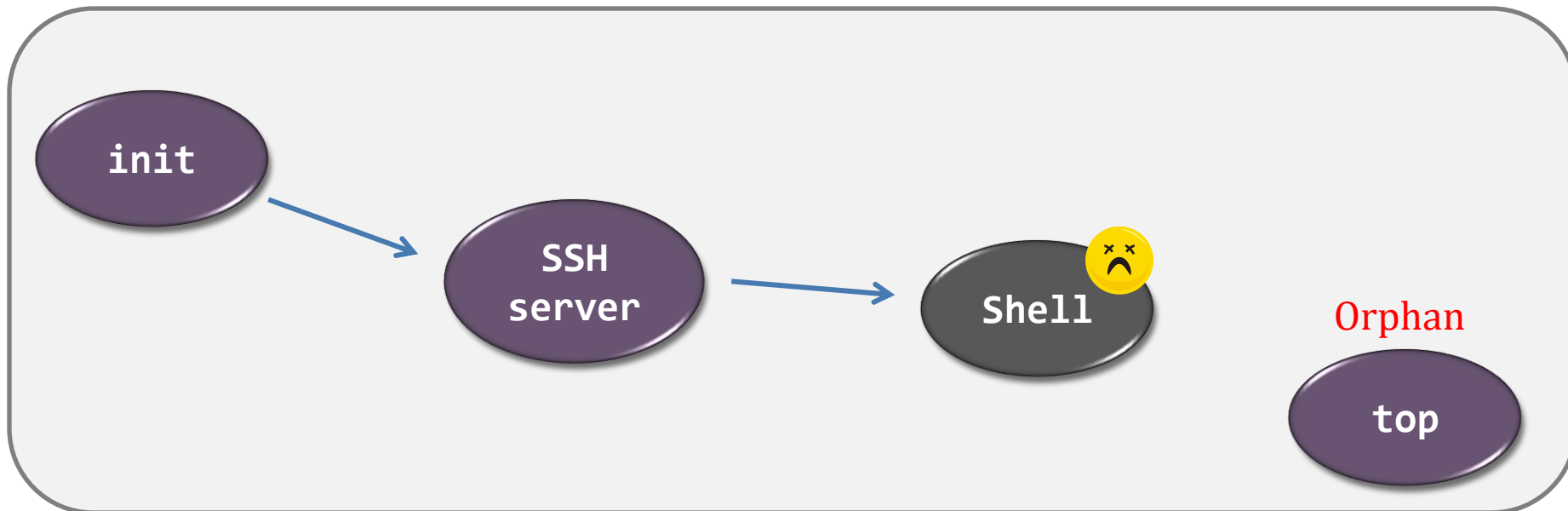  - Using **fork()** and **exec*()**.

# Process blossoming

◈ You can view the tree with the command:
  ◈ "**pstree**"; or
  ◈ "**pstree -A**" for ASCII-character-only display.



also implies the parent-child relationship.

init

fork()
& exec*()

SSH
server

fork()
& exec*()

Shell

fork()
& exec*()

top

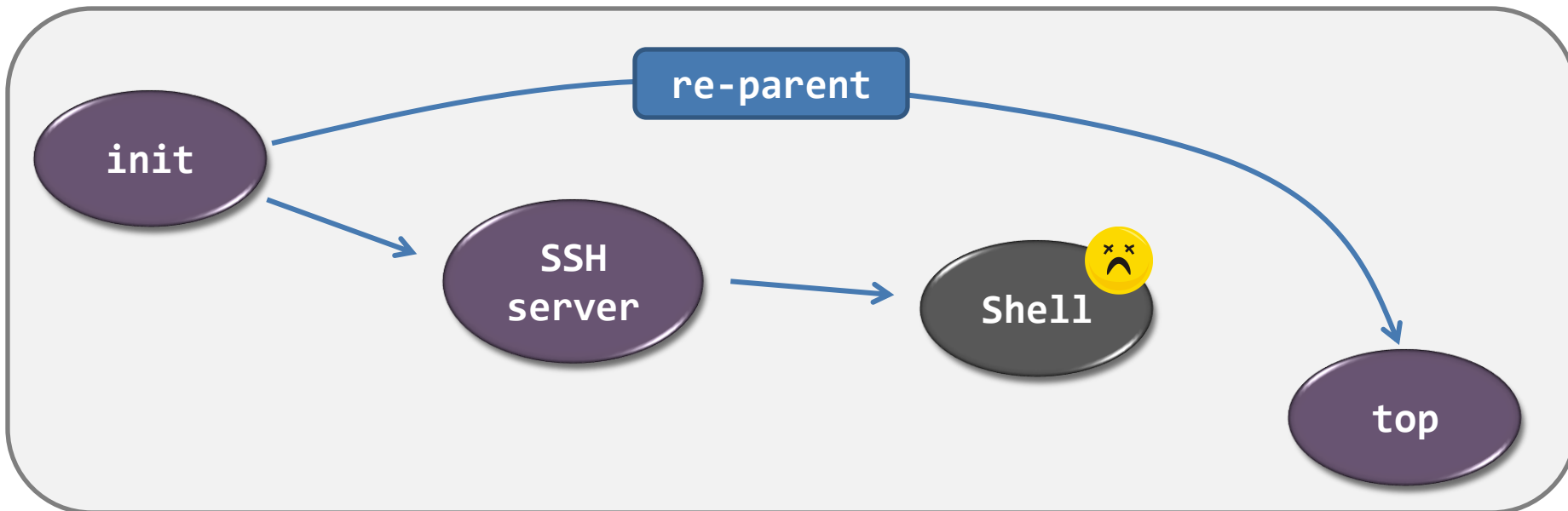# Process blossoming…with orphans?

- However, termination can happen, at any time and in any place…
  - This is no good because an orphan turns the hierarchy from a **tree** into a **forest**!
  - Plus, no one would know the termination of the orphan.

# Process blossoming…with re-parent!

- In Linux
  - The "`init`" process will become the step-mother of all orphans
  - It's called **re-parenting**

- In Windows
  - It maintains a *forest-like process hierarchy* …………



*New Linux kernels may choose someone else (e.g., the grandparent, user-level init)

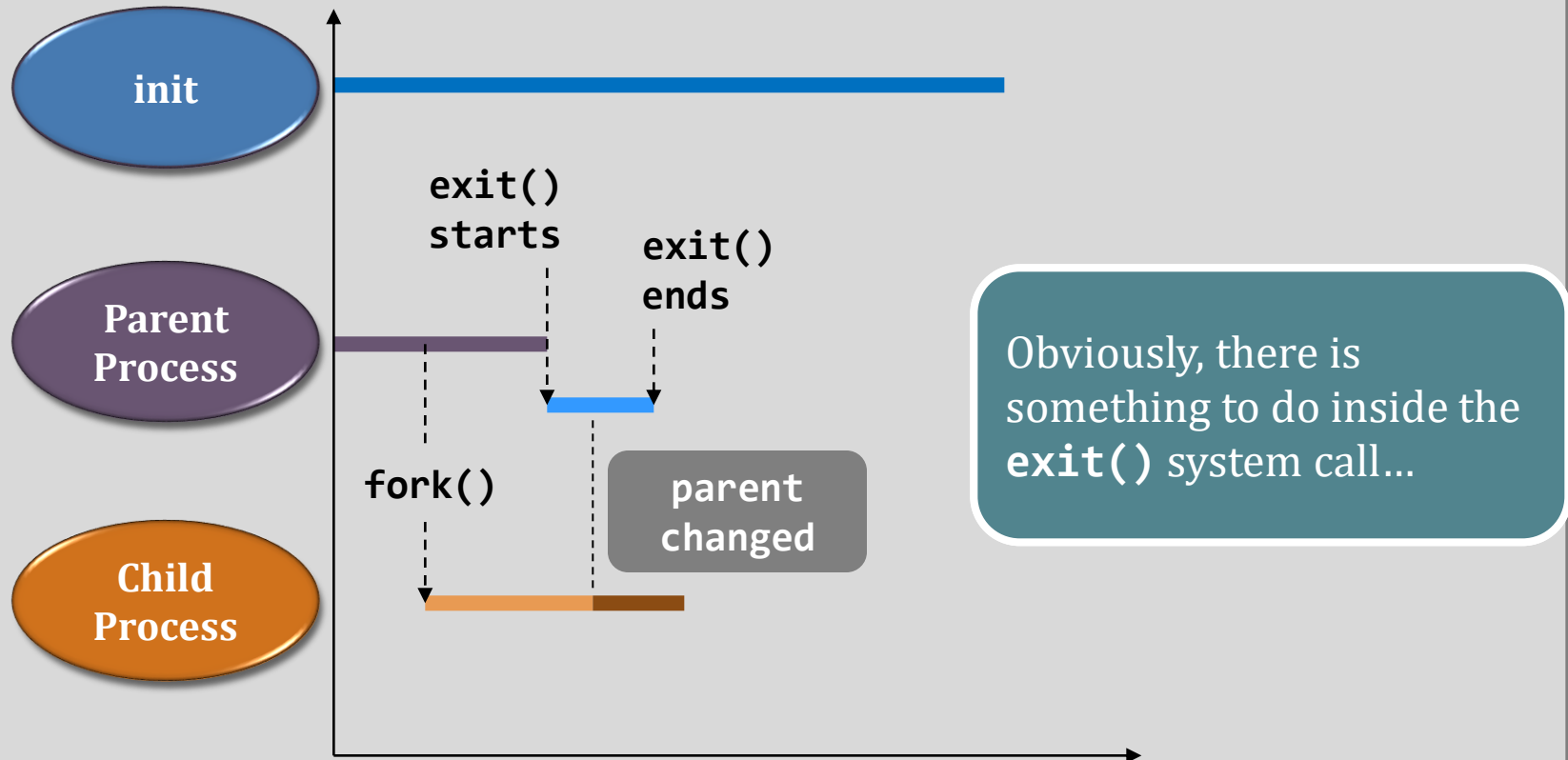# Re-parenting example

```
1   int main(void) {
2       int i;
3       if(fork() == 0) {
4           for(i = 0; i < 5; i++) {
5               printf("(%d) parent's PID = %d\n",
6                       getpid(), getppid() );
7           sleep(1);
8           }
9       }
10      else
11          sleep(1);
12      printf("(%d) bye.\n", getpid());
13  }
```

```
$ ./reparent
(1235) parent's PID = 1234
(1235) parent's PID = 1234
(1234) bye.
$ (1235) parent's PID = 1
(1235) parent's PID = 1
(1235) parent's PID = 1
(1235) bye.
$ _
```

**getppid()** is the system call that returns the parent's PID of the calling process.

# What had happened during re-parenting?

# What had happened during re-parenting?

# What had happened during re-parenting?



OS Kernel

This guy invoked **exit()**.

Process 1 ... Process 1234 Process

For each of the children of process 1234, the **exit() system call** changes its parent pointer to the "**init** process".

PID = 1
list of children
parent pointer

PID = 1234
list of children
parent pointer

PID = 1235
list of children
parent pointer

The "**init** process" accepts its new child by adding the new child into the "**list of children**".

# Background jobs

◈ The re-parenting operation enables something called **background jobs** in Linux

   ◈ It allows a process runs **without a parent terminal/shell**

**Back to home**

```
$ ./infinite_loop &
$ exit

[ The shell is gone ]
```

```
$ ps –C infinite_loop
 PID  TTY
1234  ... ./infinite_loop
$ _
```

Will see more in detail soon

# Process lifecycle

**The birth of a process.**

Except the first process "**init**", every process is created using **fork()**.

**Just fork()-ed**

**Zombie** (or terminated)

**Ready**

**Running**

**Interruptible**

**Un-interruptible**

**Blocked/Waiting**

# Process lifecycle - Ready

**Just fork()-ed**
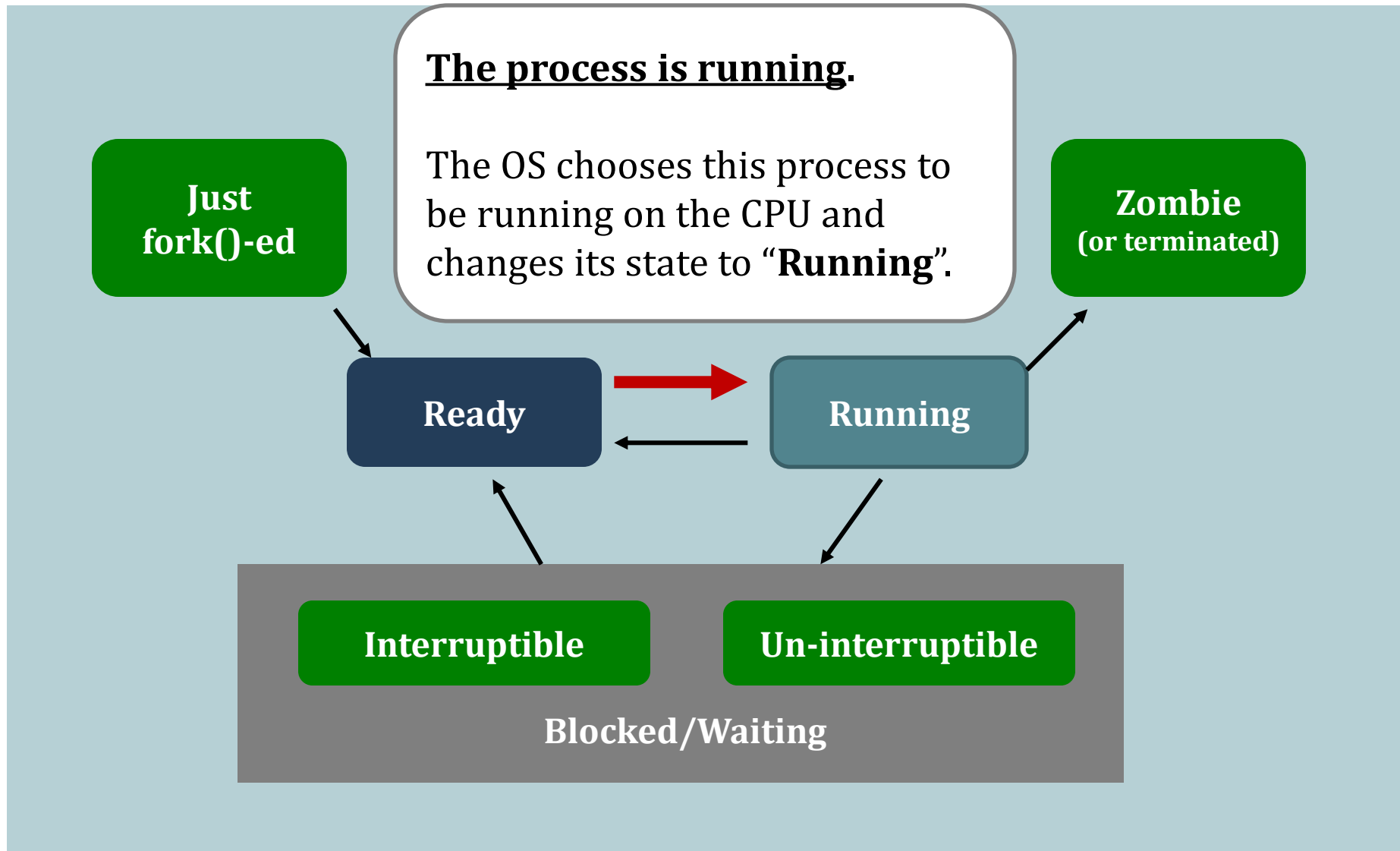
**Ready**

**Interruptible**

**Blocke**

**The process is ready**.
It means it is *ready to run* **but is not running**.

A process may become "**ready**" (*runnable*) after...
- it is just created by `fork()`;
- it has been running on the CPU for some time and the OS chooses another process to run (scheduled context switch)
- returning from blocked states.

# Process lifecycle - Running

**The process is running.**

The OS chooses this process to be running on the CPU and changes its state to "**Running**".

**Just fork()-ed**

**Zombie** (or terminated)

**Ready**

**Running**

**Interruptible**

**Un-interruptible**

**Blocked/Waiting**

# Process lifecycle - Blocking

**<u>The process is blocked</u>.**

While the process is running, it may wait for something (e.g., `getc()`, `wait()`).

**Just fork()-ed**

**Zombie** (or terminated)

**Ready**

**Running**

**Interruptible**

**Un-interruptible**

**Blocked/Waiting**

# Process lifecycle – Interruptible wait

Example.  **Reading a file**.

Sometimes, the process has to wait for the response from the device and, therefore, it is **blocked**
- this blocking state is **interruptible**
    - E.g., "**Ctrl + C**" can get the process out of the waiting state (but goes to termination state instead).

Process → fgetc(...) →

| Interruptible | Un-interruptible |
|---|---|

**Blocked/Waiting**

# Process lifecycle – Un-Interruptible wait

Sometimes, a process needs to wait for a resource until it really gets what it wants
- Doesn't want to be "Ctrl-C" interruptible
- **Un-interruptible** status
    - No way to signal it to wake up unless it returns itself
    - Check online!  The only solution is …

Who set this?
- E.g., syscall call (http://man7.org/linux/man-pages/man2/delete_module.2.html)

Why set this?
- Easier programming for lazy programmer (e.g., a driver program for a DVD drive)
- The programmer "thinks" the wait is very short and robust
    - This is one the top reasons that hang your machine / process today!
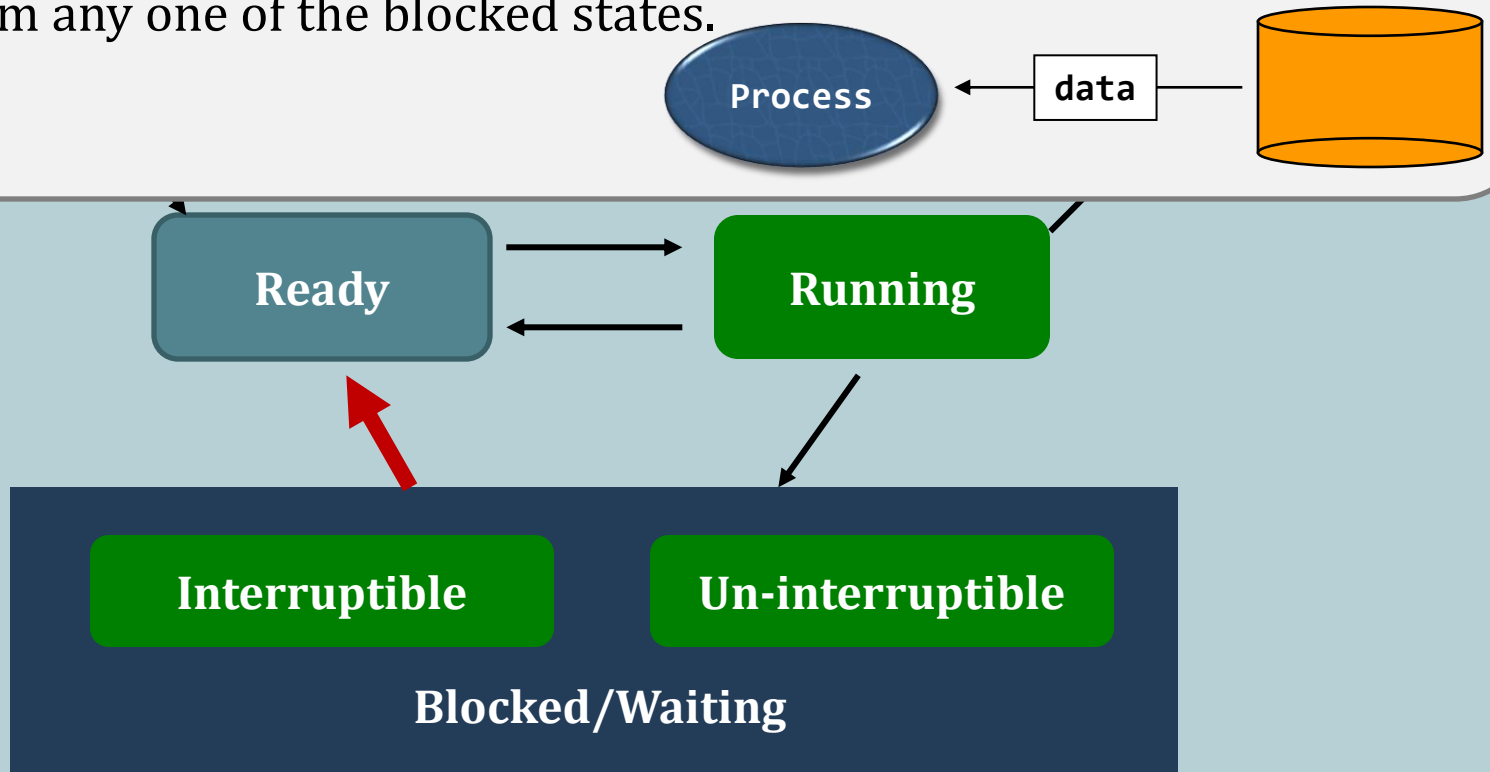- …

| Interruptible | Un-interruptible |
|---|---|

**Blocked/Waiting**

http://unix.stackexchange.com/questions/96797/what-does-the-interruptible-sleep-state-indicate

http://stackoverflow.com/questions/767551/how-to-stop-uninterruptible-process-on-linux

# Process lifecycle

**Return back to ready.**

When response arrives, the status of the process changes back to **Ready.** from any one of the blocked states.



Ready

Running

Interruptible

Un-interruptible

Blocked/Waiting

# Process lifecycle

**The process is going to die**.

The process may
- choose to terminate itself; or
- force to be terminated.

**Zombie**
**(or terminated)**

**Running**

**Interruptible**

**Un-interruptible**

**Blocking / Waiting**

# Thank You!