

# Programación en C

Alocación de memoria

# Alocación de memoria

Vimos que podemos usar pointers para manipular memoria de la manera que deseemos. También vimos cómo es que los distintos tipos de datos se guardan en memoria.

La pregunta es: ¿De donde sale esta memoria?

- ¿Cómo es que C nos da una porción de memoria para usar?
- ¿Cómo obtiene C esa memoria?
- ¿Dónde están en memoria el resto de los **procesos**?

Proceso: instancia de un programa en ejecución

# El espacio de direcciones

Cada proceso tiene acceso a un espacio de direcciones de memoria.

La cantidad máxima de memoria que podemos utilizar está limitada al tamaño de nuestro espacio de direcciones. Si no tenemos una dirección para usar, no tenemos donde guardar datos.

A la derecha tenemos un diagrama de un espacio de memoria que iremos completando con las distintas secciones que usa un proceso.

0xffffffff

0x00000000



# El espacio de direcciones

La arquitectura donde estamos trabajando determina el tamaño de nuestro espacio de direcciones.

Más específicamente, esto está dado por el tamaño de palabra del procesador. Un CPU solo puede referenciar memoria mediante una dirección que entre en una palabra.

Tamaño de palabra más común hoy en día: 64 bits.

El tamaño de los pointers está atado al tamaño de palabra.

0xffffffff

0x00000000



# El espacio de direcciones

Un pointer de 32 bits puede tomar valores entre: `0x0` y `0xffffffff`. Esto le permite acceder a 4GB de memoria.

Un pointer de 64 bits entre: `0x0` y `0xffffffffffffffff`. Esto se traduce a 16 exabytes. 1 exabyte =  $10^9$  GB.

En el diagrama asumimos 32 bits para tener una notación más compacta.

`0xffffffff`



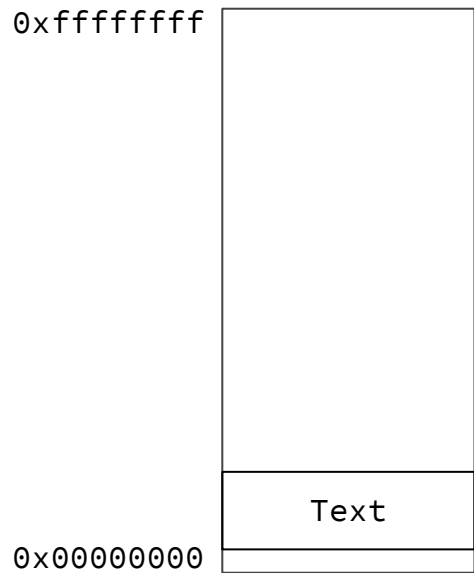
`0x00000000`

# Código

Primero que nada, el proceso necesita tener acceso al programa que va a ejecutar.

El código del programa (en binario, ya compilado), junto con todas las librerías que use, se aloca en la sección llamada **Text**.

El resto de la memoria se usa para datos. Dependiendo cómo se aloquen en el programa estos datos, van a terminar en distintas secciones.



# Tiempo de vida (lifetime)

C99: “El tiempo de vida de un objeto es el período de tiempo en la ejecución del programa durante el cual se garantiza que va a tener almacenamiento.

Un objeto existe, tiene una dirección constante y mantiene su último valor guardado durante su tiempo de vida.

Si se refiere a un objeto fuera de su tiempo de vida, el comportamiento es indefinido”.

El tiempo de vida de un objeto está determinado por la duración de su almacenamiento (storage duration).

# Duración de almacenamiento

El estándar define tres duraciones de almacenamiento:

- Estática
- Automática
- Alocada

Para evitar confusiones nos referiremos a la duración “alocada” como “dinámica”.

Este es el nombre que usa el estándar de C++ y es más claro. Al fin y al cabo, todo almacenamiento debe ser **alocado** en memoria.

De las 3 duraciones surgen 3 tipos de alocaiones.



# Alocación estática

# Alocación estática

C99: “El tiempo de vida del objeto es la ejecución entera del programa. Su valor es inicializado solo una vez, previo al inicio del programa”.

1. El objeto siempre existe
2. El objeto siempre va a ser inicializado.
3. Todo esto se garantiza antes de que corra la función **main**.

¿Qué objetos son alocados estáticamente?

- Variables declaradas fuera de una función.
- Variables dentro de una función declaradas como **static**.

# Alocación estática: inicialización

C99: “Un objeto con alocación estática sólo puede ser inicializado con una constante.

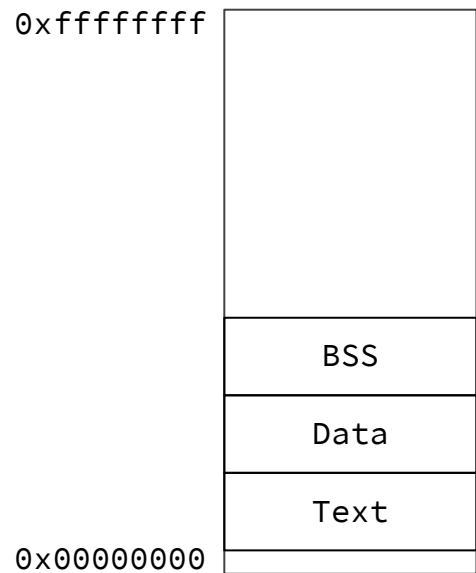
Si no es inicializado explícitamente:

- Si es un pointer, es inicializado a NULL.
- Si es un tipo aritmético, es inicializado a 0.
- Si es un array/struct, todos los miembros son inicializados siguiendo estas reglas.
- Si es una union, el primer miembro es inicializado siguiendo estas reglas.”

# Alocación estática

Los objetos alocados estáticamente tienen sus propias secciones de memoria. Estas son de tamaño fijo ya que el espacio que ocupan es conocido en tiempo de compilación.

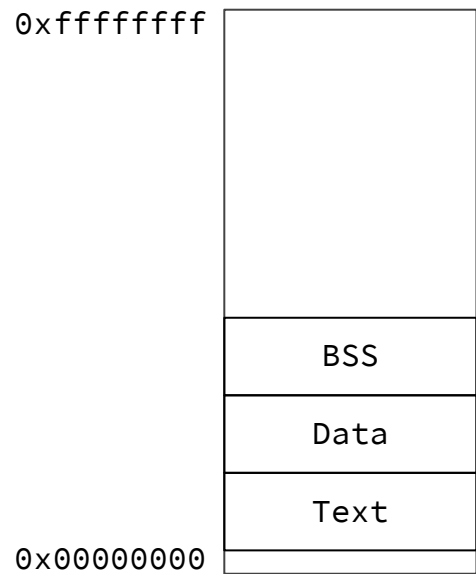
- BSS: para las variables no inicializadas.
- Data: para las variables inicializadas.



# Alocación estática

Cuando el proceso arranca la sección Data es copiada tal cual se encuentra en el ejecutable directo a memoria.

Cómo el BSS por definición está compuesto puramente de ceros, es generado en el momento con el tamaño especificado en el ejecutable.



# Alocación estática

Cómo la estructura de las dos secciones de memoria están fijas en el ejecutable, la memoria en esas secciones va a ser idéntica para todos los procesos que corran el mismo programa.

Por lo tanto, todos los objetos que están alocados estáticamente **siempre** van a estar en las mismas posiciones **relativas** entre sí.

La posición **absoluta** depende de varios factores y puede cambiar entre diferentes ejecuciones de un mismo programa, pero es algo que no vamos a cubrir.

Temas relacionados: Position Independent Code, ASLR.

# Alocación estática

Si yo declaro dos variables:

```
char a;
```

```
int b;
```

Y el compilador decide que a se aloca antes que b.



$\&a < \&b$  en **todos** los procesos.

# Alocación estática

A diferencia de un struct donde el orden de declaración de las variables importa, el compilador puede ordenar las variables para minimizar padding o hacer otras optimizaciones.

**char** a;

**int** b;

**char** c;

**char** d;

**char** e;





# Problema: múltiples procesos

¿Qué pasa si corremos dos procesos a partir del mismo programa de manera simultánea?

Sabemos que las secciones donde se encuentran los objetos alocados estáticamente van a ser idénticas y posiblemente estén en las mismas direcciones.

¿Qué pasa cuando ambos procesos acceden a la misma variable y por lo tanto a la misma **dirección de memoria**?

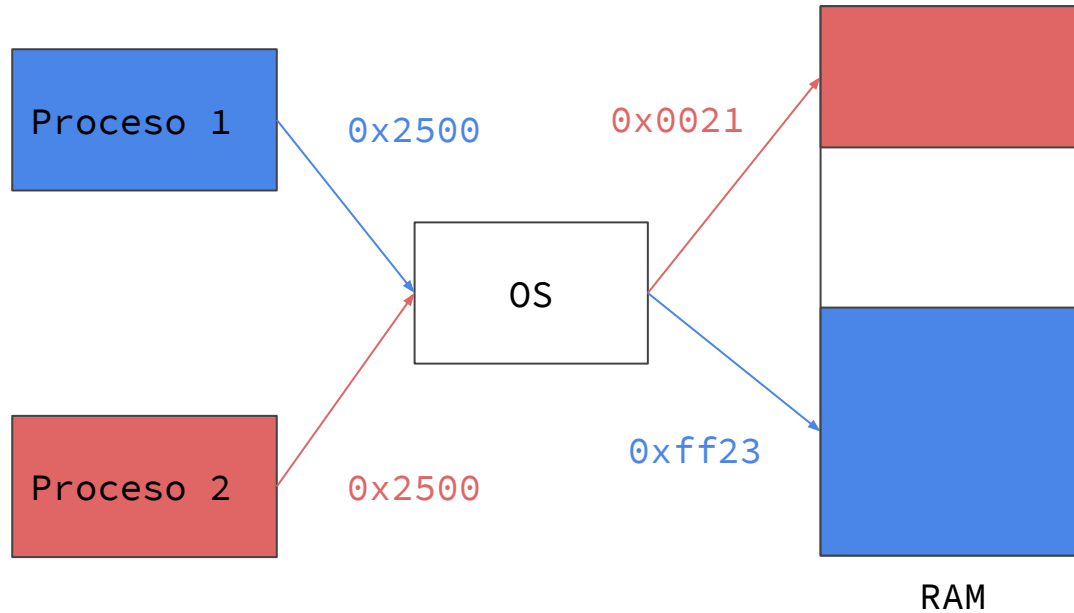
# Memoria virtual

Una de las funcionalidades de un Sistema Operativo es la memoria virtual.

Cuando un proceso inicia, el SO le brinda un espacio de direcciones completamente privado, de manera que ningún proceso puede interactuar con el espacio de memoria de otro.

Cada proceso puede utilizar esta memoria sin preocuparse por el resto de las cosas corriendo en el sistema. El SO se encarga de traducir las direcciones virtuales de cada proceso a direcciones de memoria física.

# Memoria virtual



# Memoria virtual

Sin memoria virtual, trabajar en un lenguaje como C con acceso a pointers es propenso a errores y además muy inseguro.

Podríamos:

- Leer memoria de otro proceso. Ej: contraseñas guardadas en el navegador.
- Crashear otro proceso, incluido el sistema operativo entero.
- Sobrepasar cualquier tipo de “lock screen”.
- Y muchas cosas más...

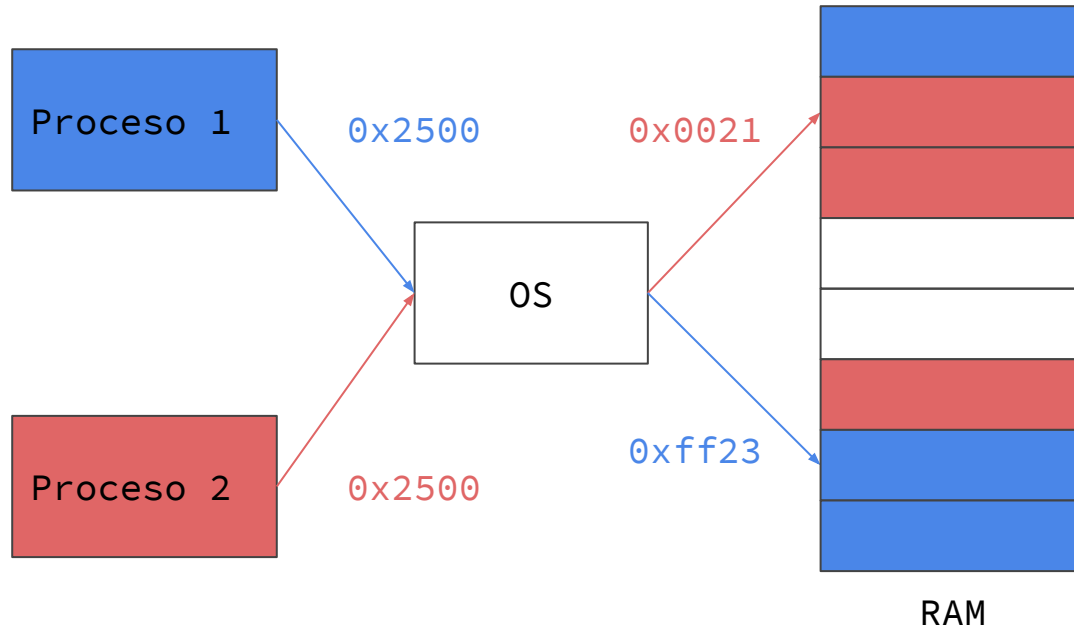
# Alocación de memoria virtual

El SO no le da todo el espacio de memoria a todos los procesos. Cuando un proceso necesita memoria le debe pedir al SO.

El SO asigna memoria en páginas. Cada página tiene un tamaño fijo, el más común es 4kb (Linux, Windows, MacOS). Del lado del proceso una página que le fue asignada por el SO se dice que está **mapeada**.

Cuando el proceso arranca, va a solicitar suficiente memoria como para alojar todas las variables estáticas más el código del programa.

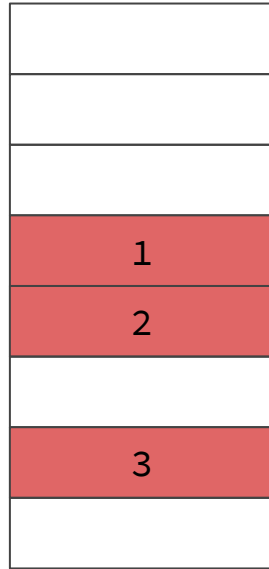
# Memoria virtual



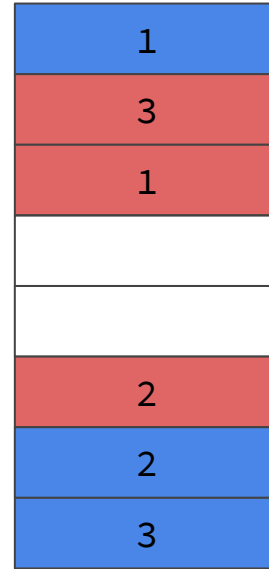
# Memoria virtual



Memoria  
Proceso 1



Memoria  
Proceso 2



RAM

Página en blanco == No mapeada

# Permisos

Cada página puede tener distintos tipos de permisos de acceso, dependiendo lo que requiera el proceso.

Los permisos son: read, write y execute.

**Segmentation Fault** es una señal generada por el SO. Ocurre cuando el proceso intenta realizar una operación en una página que no tiene los suficientes permisos. Esto incluye cualquier tipo de operación realizada en una página no mapeada.

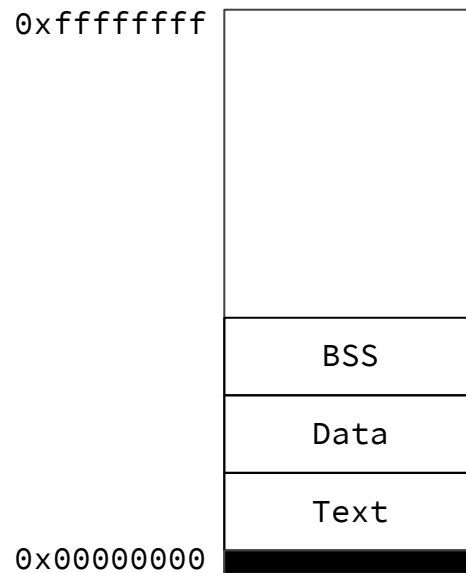


# NULL

En la **mayoría de los casos**, el compilador va a elegir la dirección 0 como el null pointer.

Para facilitar la detección de la desreferenciación de un null pointer, la primera página de memoria nunca es mapeada por el proceso, garantizando un Segmentation Fault.

Esto sigue siendo **undefined behavior** y puede no ocurrir en ciertos casos, particularmente en sistemas embebidos, donde la memoria es escasa.

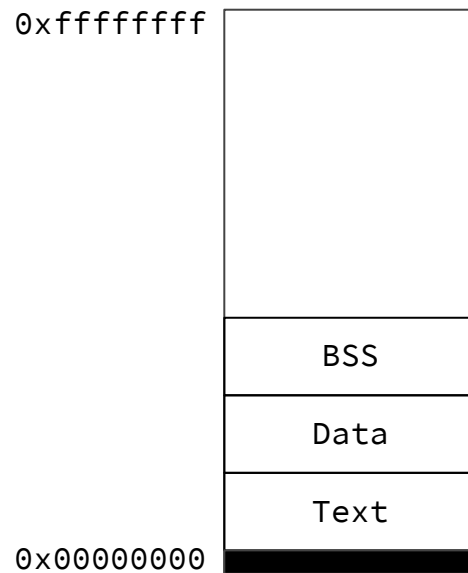


# Secciones != páginas

Las distintas secciones de memoria que crea el proceso son puramente lógicas y pueden ocupar más de una página.

Las secciones BSS y Data también pueden encontrarse en la misma página ya que tienen los mismos permisos.

Sin embargo, Text siempre debe estar en páginas separadas de BSS y Data ya que requiere **execute**.

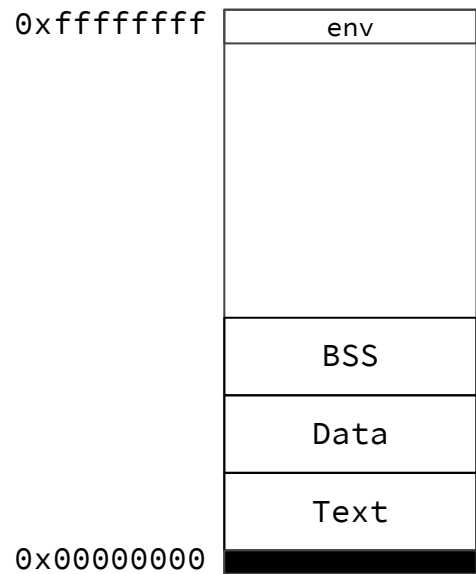


# Entorno

En las direcciones de memoria más altas se aloca espacio para el entorno del programa.

Esto incluye:

- Variables de entorno
- Argumentos de la línea de comandos



# Alocación automática

# Alocación automática

C99: “El tiempo de vida del objeto va desde la entrada al bloque con el que está asociado hasta que la ejecución de ese bloque termine de cualquier manera (Ingresar en un sub-bloque o llamar una función suspende, pero no termina, la ejecución del bloque actual).

Si se ingresa al bloque recursivamente, una nueva instancia del objeto es creado cada vez.

El valor inicial del objeto es indefinido, a menos que sea inicializado.”

Bloque: todo el código que esté encerrado entre { }.

# Alocación automática

¿Qué objetos son alocados estáticamente?

Toda variable declarada dentro de un bloque (función, loop, etc).

Siempre y cuando no sean declaradas **static** (o **extern**).

No son inicializados automáticamente

Leer su valor antes que sea inicializado es undefined behavior.

Solo existen **dentro el bloque**

# Alocación automática

En alocación estática, sabemos en tiempo de compilación exactamente cuántas variables vamos a tener junto con su tamaño. En cambio, en alocación automática vamos a necesitar más o menos memoria dependiendo que pase en runtime.

Cada vez que entramos a un bloque nuevo necesitamos alocar memoria para sus variables.

- Llamar a una función
- Entrar en un loop
- Entrar en un if, o en un else

# Alocación automática

En vez de una sección de memoria fija, las variables automáticas deben vivir en una sección que sea también variable en runtime.

Esta sección se crea al arrancar el programa y crece según es necesario para acomodar nuevas alocaiones. Si la sección pasa a sobrepasar el tamaño de la página, otra página debe ser solicitada al sistema operativo para que pueda continuar creciendo.



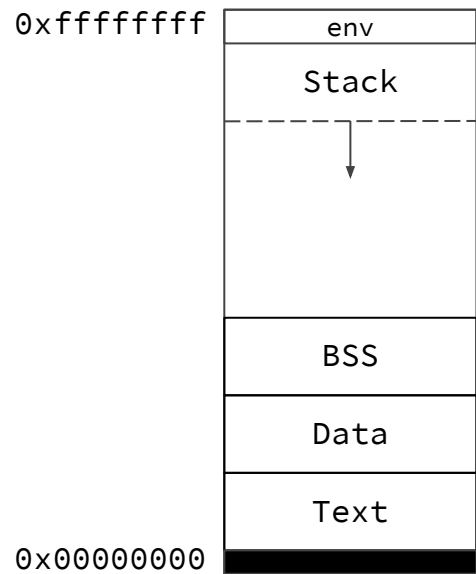
# El Stack

Probablemente conocen StackOverflow, ya sea la página o la excepción de Java.

Estos hacen alusión a este stack, usado por la gran mayoría de los programas.

El propósito del stack es guardar valores temporales para facilitar la ejecución del programa:

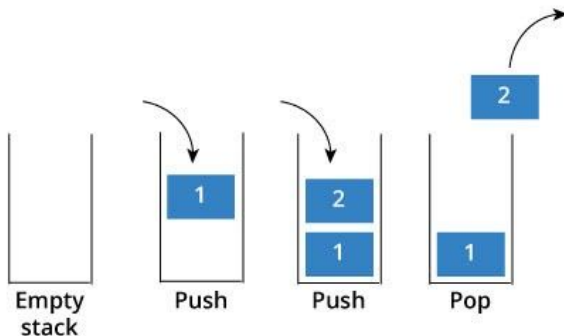
- Alocaciones automáticas
- Dónde volver cuando ejecutamos **return**



# ¿Qué es un stack?

Un stack es una estructura de datos que funciona como una colección de objetos, con dos operaciones:

- Push: agrega un objeto al stack
- Pop: remueve el último objeto agregado al stack

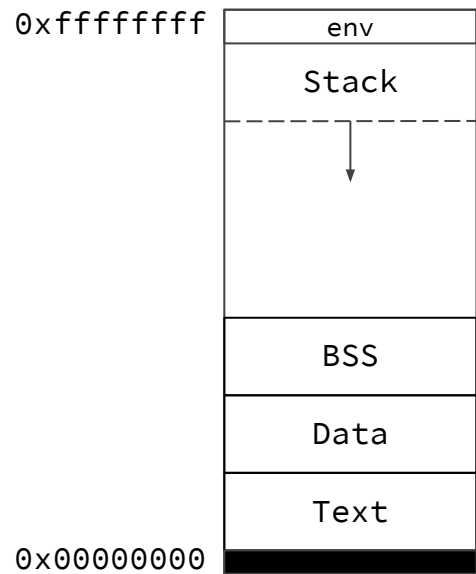


# El Stack

“El Stack” hablando de la sección de memoria, implementa esta misma funcionalidad.

El inicio o “bottom” del stack se encuentra en las direcciones de memoria altas. Los elementos se agregan de manera que crece hacia las direcciones de memoria bajas.

Se podría decir que está “invertido”.

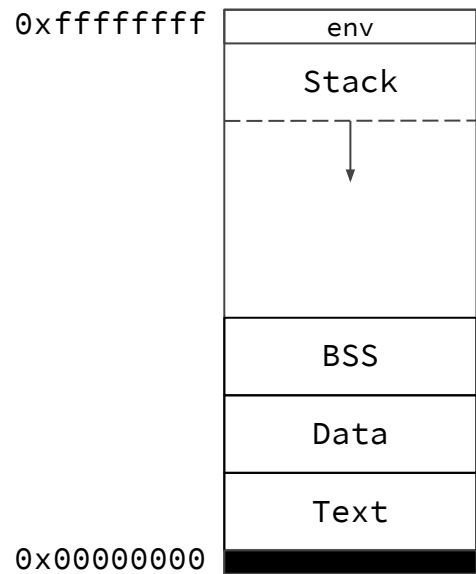


# El Stack

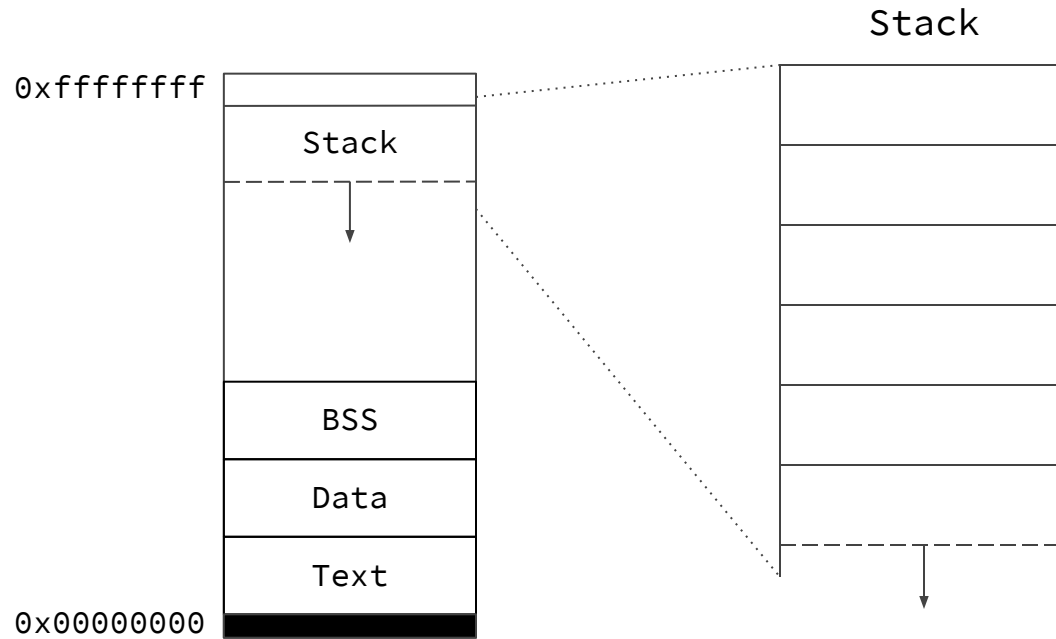
A diferencia de las otras secciones, el stack va a crecer durante la ejecución del programa.

Al iniciar el proceso se alocan una cierta cantidad de páginas. Conforme el stack crece el SO automáticamente le aloca nuevas páginas. Esto se repite hasta que se llega al tamaño máximo del stack.

Lo que en Java se notifica con un `StackOverflow`, en C probablemente se vea como un `Segmentation Fault` debido a utilizar una página no mapeada.



# El Stack



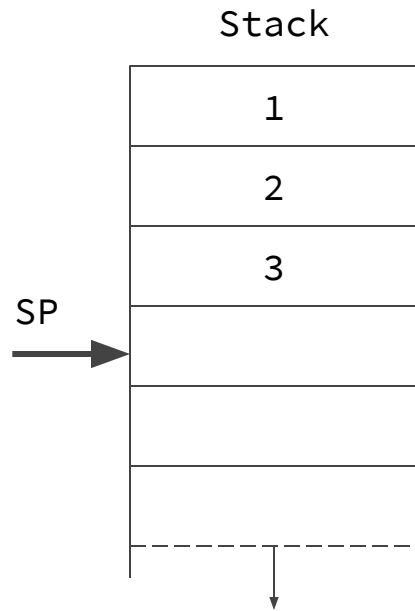
# El Stack

Se representa como una sección de memoria contigua.

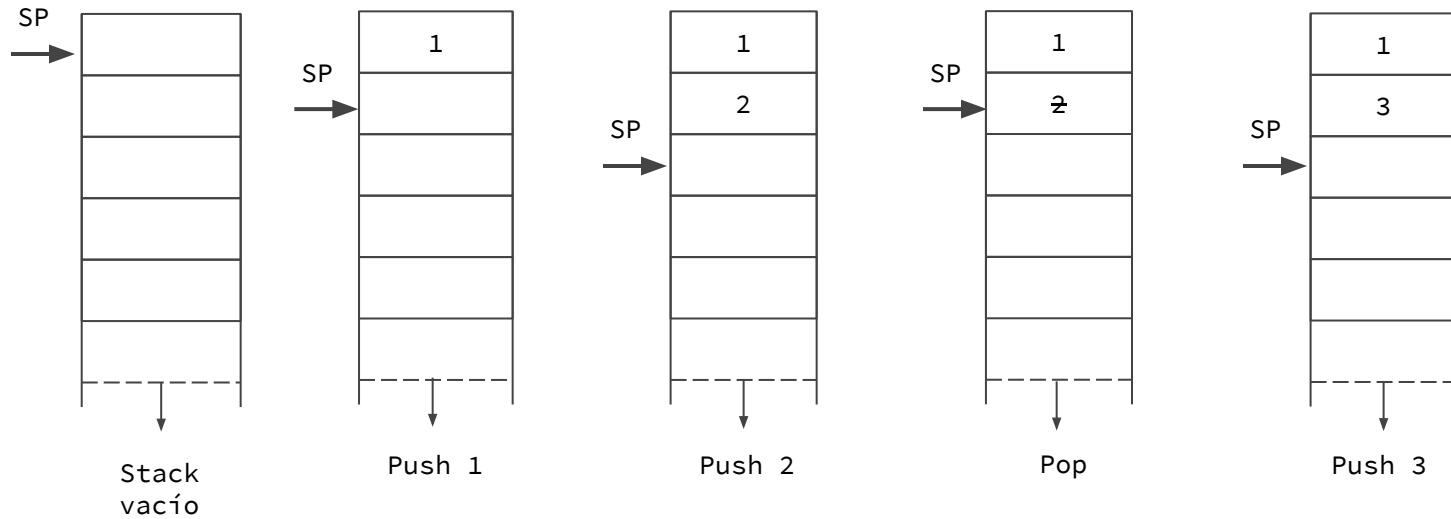
El bottom del stack es fijo, por lo que no necesitamos actualizar su posición. El top del stack, sin embargo, cambia dependiendo cuantos elementos tiene.

Stack Pointer (SP): pointer al top del stack

Conforme vamos agregando elementos, avanzamos el SP. Para remover un elemento, simplemente hacemos retroceder el SP.



# El Stack



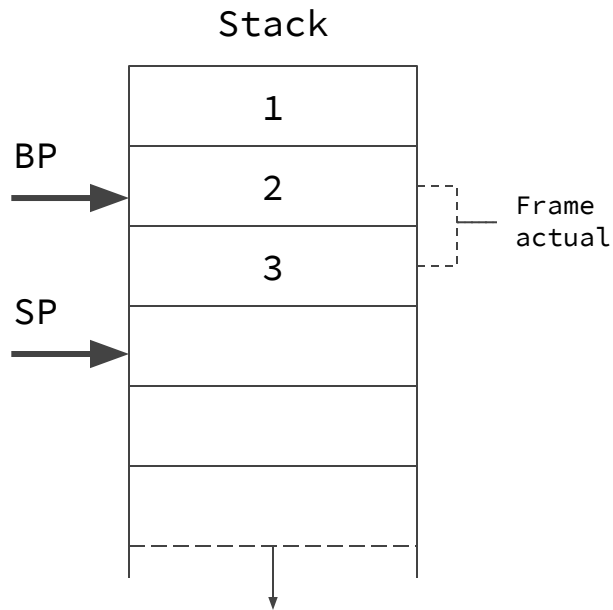
El SP apunta al  
primer lugar vacío.

# Function frame

Un function frame es la parte del stack que pertenece a una función. Cada función que se llama crea un nuevo function frame.

Para marcar el function frame de la función que se está ejecutando en este momento, utilizamos otro pointer.

Base pointer (BP): marca el inicio del function frame actual.





# Function frame

Las variables alocadas automáticamente son alocadas dentro del function frame de cada función. Cada vez que llamamos a una función se va a crear un nuevo frame con el espacio suficiente para guardar todas sus variables.

Esta alocaión se hace simplemente avanzando el Stack Pointer la cantidad de bytes que sean necesarios. Esto se hace decrementando el valor del pointer. Si quiero alocar 8 bytes:  $SP = SP - 8$  (pointer arithmetic como **char\***).

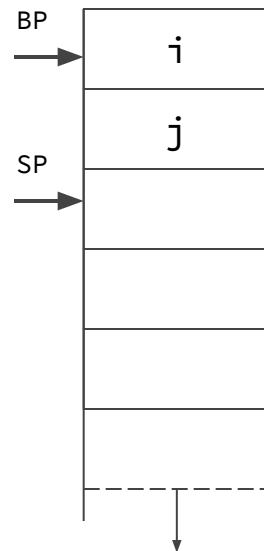
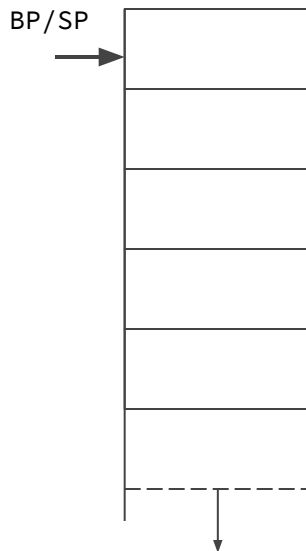
Para referenciar una variable, el compilador usa un offset desde el Base Pointer.

# Function frame

```
void function() {  
    int i;  
    int j;  
}
```

Para ser más acotados en la slide:

- Cada espacio en el stack son 4 bytes.
- Seguimos usando 32 bits. Los pointer ocupan 4 bytes.



$SP = SP - 8$

El compilador asigna:

$\&i = BP$

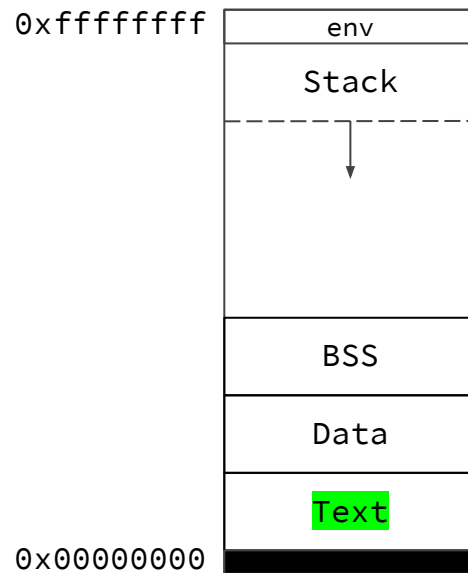
$\&j = BP - 4$

# Instruction pointer

Recordamos que la sección Text contiene el código que está ejecutando el programa.

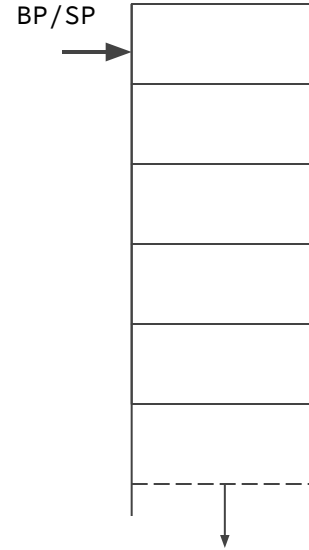
Además del Stack Pointer y el Base Pointer para marcar el frame de función actual, también debemos mantener un pointer a la instrucción que estamos ejecutando. Por instrucción nos referimos a código binario, pero podemos verlo como la línea del source code ahora.

Instruction Pointer (IP): marca la instrucción que se está ejecutando.



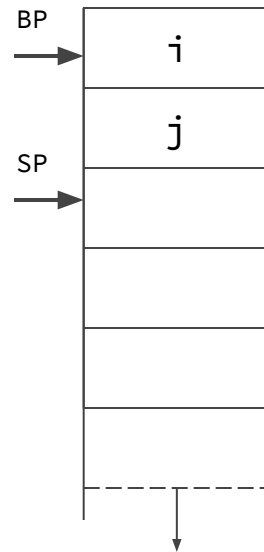
# Ejecución

```
IP → 00 void function() {  
      01     int i = get_five();  
      02     int j = get_five();  
      03 }  
      04  
      05 int get_five() {  
      06     int five = 5;  
      07     return five;  
      08 }
```



# Ejecución

```
IP → 00 void function() {  
      01     int i = get_five();  
      02     int j = get_five();  
      03 }  
      04  
      05 int get_five() {  
      06     int five = 5;  
      07     return five;  
      08 }
```



IP++

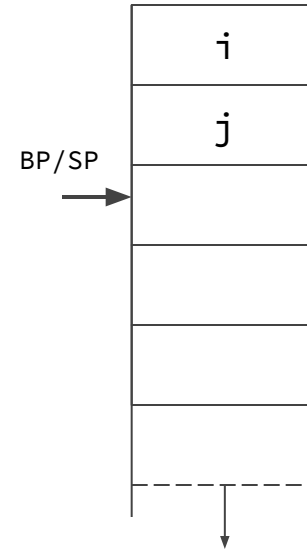
SP = SP - 8

IP++: próxima  
instrucción.

El CPU lo ejecuta  
automáticamente.

# Ejecución

```
00 void function() {  
01     int i = get_five();  
02     int j = get_five();  
03 }  
04  
IP → 05 int get_five() {  
06     int five = 5;  
07     return five;  
08 }
```



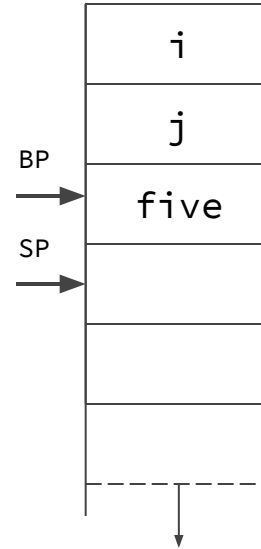
IP = &get\_five

BP = SP

# Ejecución

```
00 void function() {  
01     int i = get_five();  
02     int j = get_five();  
03 }  
04  
05 int get_five() {  
06     int five = 5;  
07     return five;  
08 }
```

IP  
→



IP++

SP = SP - 4

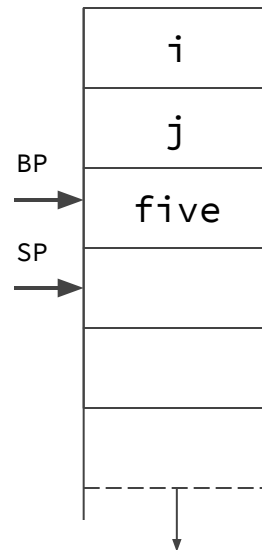
# Ejecución

```
00 void function() {  
01     int i = get_five();  
02     int j = get_five();  
03 }  
04  
05 int get_five() {  
06     int five = 5;  
07     return five;  
08 }
```

IP



???



IP = ??

BP = ??

SP = ??

???



# Ejecución

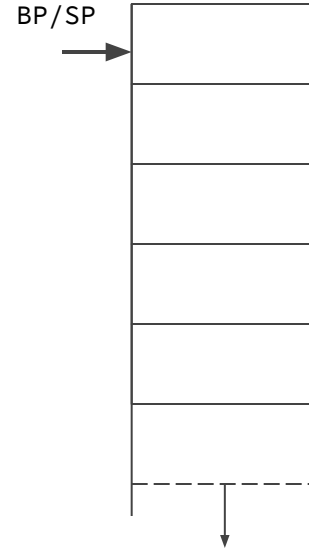
Después del **return** necesitamos una manera de poder restaurar el function frame anterior y de ejecutar la próxima instrucción.

La única manera que podemos hacer esto es guardando los valores que necesitamos antes de sobrescribirlos.

¿Dónde los guardamos? ¡En el stack!

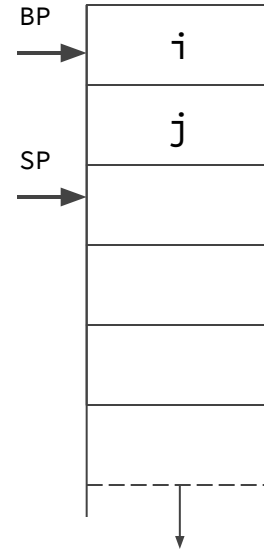
# Ejecución

```
IP → 00 void function() {  
      01     int i = get_five();  
      02     int j = get_five();  
      03 }  
      04  
      05 int get_five() {  
      06     int five = 5;  
      07     return five;  
      08 }
```



# Ejecución

```
IP → 00 void function() {  
      01     int i = get_five();  
      02     int j = get_five();  
      03 }  
      04  
      05 int get_five() {  
      06     int five = 5;  
      07     return five;  
      08 }
```

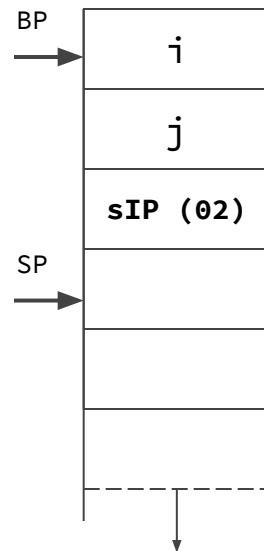


IP++

SP = SP - 8

# Ejecución

```
IP → 00 void function() {  
      01     int i = get_five();  
      02     int j = get_five();  
      03 }  
      04  
      05 int get_five() {  
      06     int five = 5;  
      07     return five;  
      08 }
```



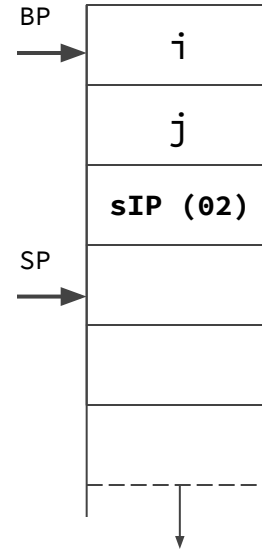
Push IP + 1 (SP -= 4)

i, j son solo  
alocaciones.

Distinguimos **sIP**  
con negrita para  
marcar que es un  
valor que  
pusheamos.

# Ejecución

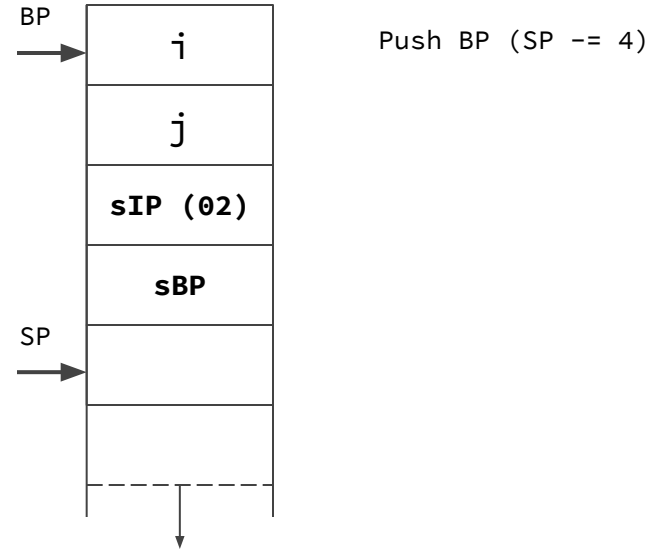
```
00 void function() {  
01     int i = get_five();  
02     int j = get_five();  
03 }  
04  
IP → 05 int get_five() {  
06     int five = 5;  
07     return five;  
08 }
```



IP = &get\_five

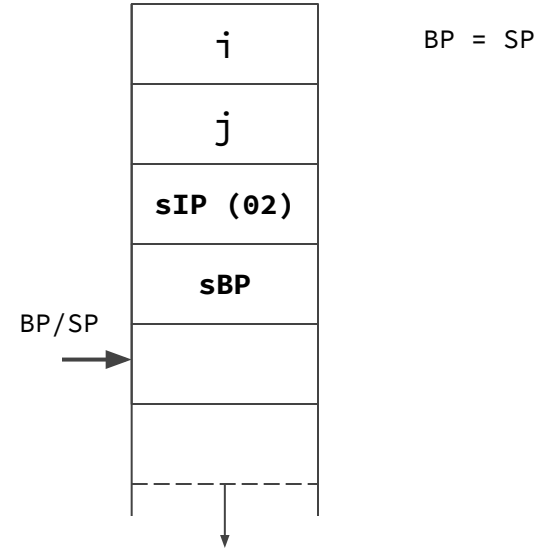
# Ejecución

```
00 void function() {  
01     int i = get_five();  
02     int j = get_five();  
03 }  
04  
IP → 05 int get_five() {  
06     int five = 5;  
07     return five;  
08 }
```



# Ejecución

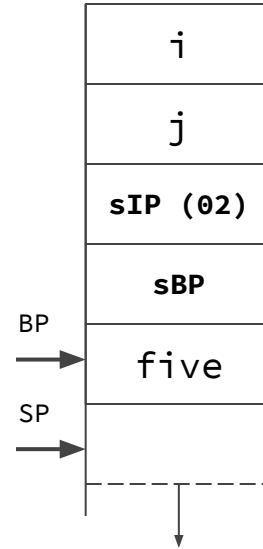
```
00 void function() {  
01     int i = get_five();  
02     int j = get_five();  
03 }  
04  
IP → 05 int get_five() {  
06     int five = 5;  
07     return five;  
08 }
```



# Ejecución

```
00 void function() {  
01     int i = get_five();  
02     int j = get_five();  
03 }  
04  
05 int get_five() {  
06     int five = 5;  
07     return five;  
08 }
```

IP  
→



IP++

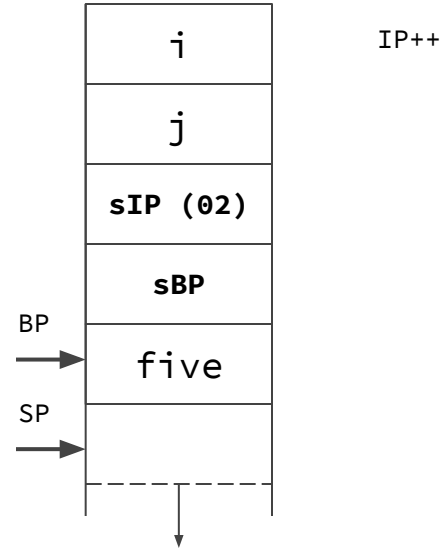
SP = SP - 4



# Ejecución

```
00 void function() {  
01     int i = get_five();  
02     int j = get_five();  
03 }  
04  
05 int get_five() {  
06     int five = 5;  
07     return five;  
08 }
```

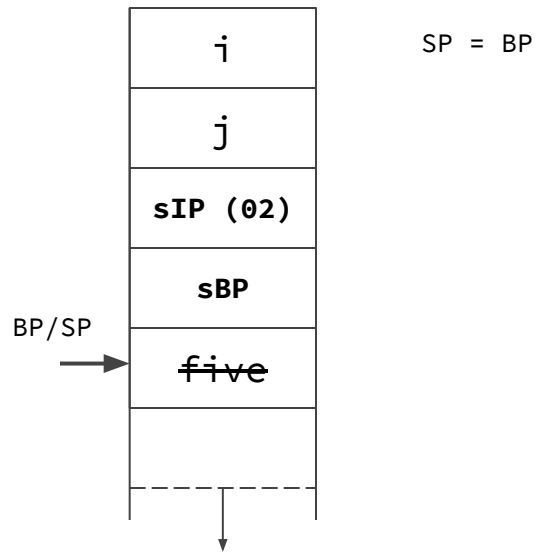
IP  
→



# Ejecución

```
00 void function() {  
01     int i = get_five();  
02     int j = get_five();  
03 }  
04  
05 int get_five() {  
06     int five = 5;  
07     return five;  
08 }
```

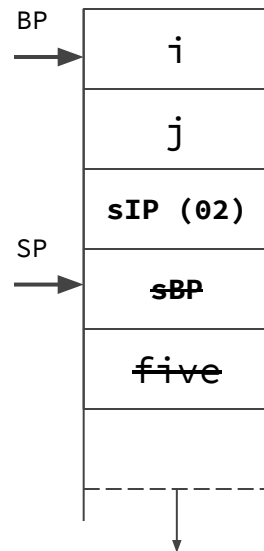
IP  
→



# Ejecución

```
00 void function() {  
01     int i = get_five();  
02     int j = get_five();  
03 }  
04  
05 int get_five() {  
06     int five = 5;  
07     return five;  
08 }
```

IP  
→

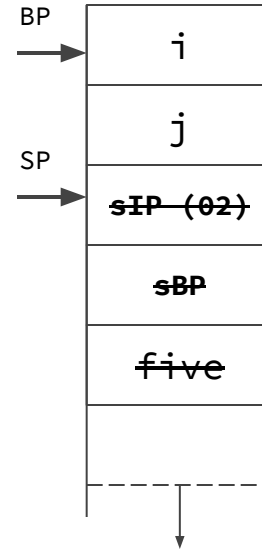


BP = Pop (SP += 4)

Pop retorna sBP

# Ejecución

```
00 void function() {  
01     int i = get_five();  
IP → 02     int j = get_five();  
03 }  
04  
05 int get_five() {  
06     int five = 5;  
07     return five;  
08 }
```



IP = Pop (SP += 4)

Pop retorna sIP

# Otros bloques

Los demás bloques que no son funciones no tienen sus propios frames.

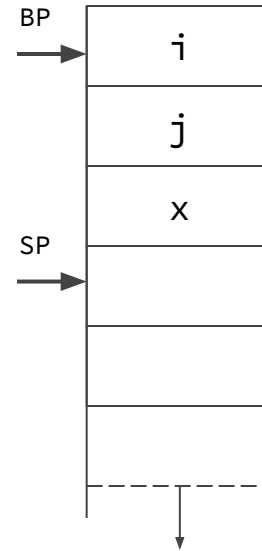
Es posible que el compilador agrande el function frame actual al entrar a un bloque y cuando el programa sale del bloque achique el frame para desocupar espacio en el stack.

También puede ser que se aloquen todas las variables al inicio de la función, sin importar a qué bloque pertenecen.

El compilador puede tomar la decisión que quiera, ya que referenciar un objeto fuera de su lifetime es **undefined behavior**.

# Otros bloques (Opción 1)

```
      00 void function() {  
IP → 01     int i;  
      02     int j;  
      03     {  
      04         int x;  
      05     }  
      06 }
```

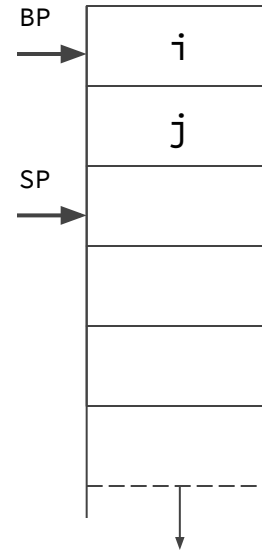


x ya está alocada  
fuera de su bloque

SP -= 12

## Otros bloques (Opción 2)

```
      00 void function() {  
IP → 01     int i;  
      02     int j;  
      03     {  
      04         int x;  
      05     }  
      06 }
```

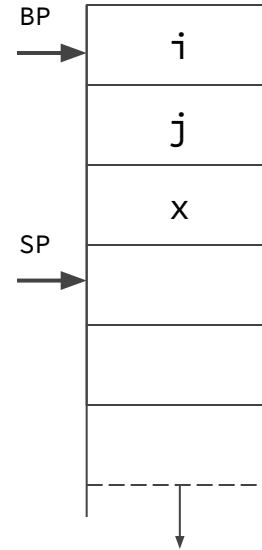


Sólo i, j están  
alocadas.

SP -= 8

## Otros bloques (Opción 2)

```
00 void function() {  
01     int i;  
02     int j;  
03     {  
IP → 04         int x;  
05     }  
06 }
```



x se aloca al  
entrar a su bloque  
  
SP -= 4



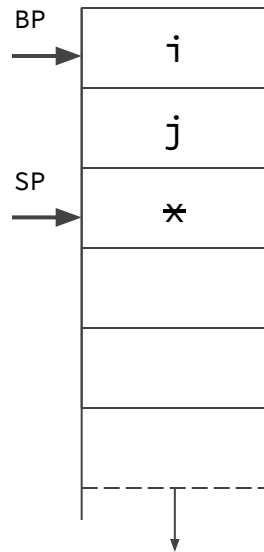
## Otros bloques (Opción 2)

```

00 void function() {
01     int i;
02     int j;
03     {
04         int x;
05     }
06 }

```

IP →



Al salir del bloque, x se desaloca.

Puedo usar ese espacio para otro bloque.

SP += 4

# Arrays en el stack

Los miembros de un array se alocan contiguamente en memoria. Conforme avanzamos en el array la dirección donde está alocao el miembro debe **subir**.

La alocaón del stack va desde direcciones altas a direcciones bajas. Sin embargo, la lectura/escritura de los datos va desde direcciones bajas a altas.

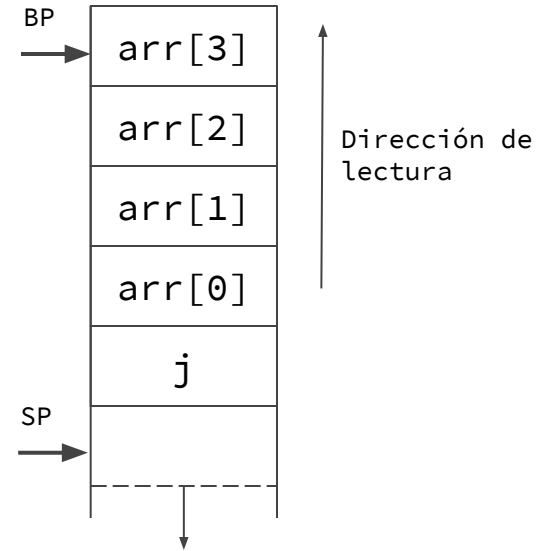
El primer elemento del array va a estar cerca del top del stack.

Si leemos/escribimos más allá de los límites del array (overflow), vamos a sobrescribir function frames anteriores.

Lo mismo se aplica a structs y unions.

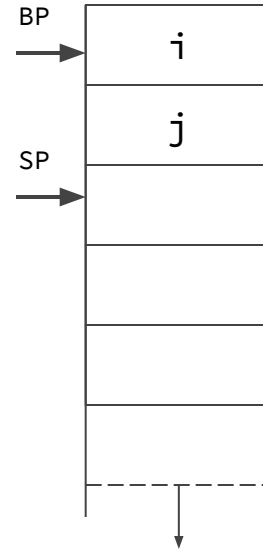
# Arrays en el stack

```
      00 void function() {  
IP → 01     int arr[4];  
      02     int j;  
      03 }
```



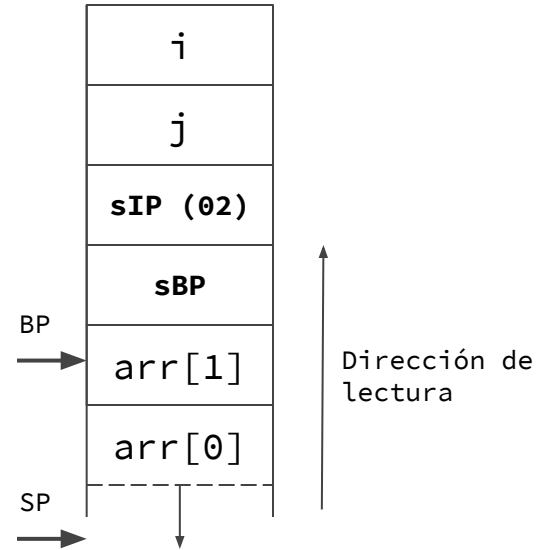
# Arrays en el stack

```
IP → 00 void function() {  
      01     int i, j;  
      02     other_function();  
      03 }  
      04  
      05 void other_function() {  
      06     int arr[2];  
      07 }  
      08
```



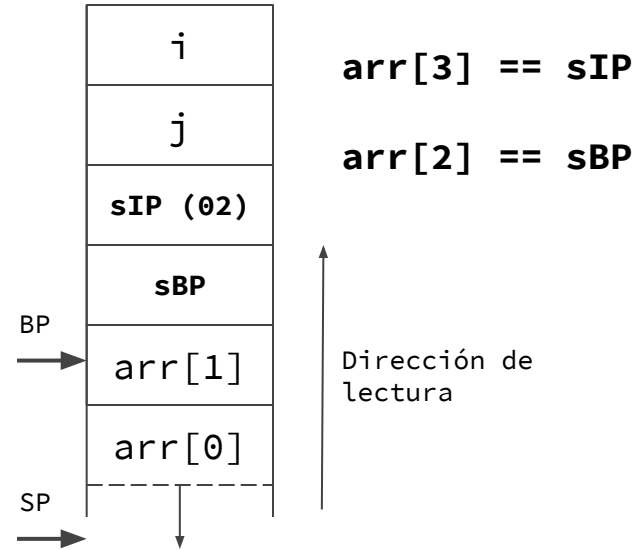
# Arrays en el stack

```
00 void function() {  
01     int i, j;  
02     other_function();  
03 }  
04  
IP → 05 void other_function() {  
06     int arr[2];  
07 }  
08
```



# Arrays en el stack

```
00 void function() {  
01     int i, j;  
02     other_function();  
03 }  
04  
IP → 05 void other_function() {  
06     int arr[2];  
07 }  
08
```



# Arrays en el stack

Escribir fuera de los límites de un array en el stack en el **mejor** de los casos causa Segmentation Fault.

En casos peores, corrupción de memoria.

En el peor de los casos, le da al usuario control total de la computadora.

Existen **mitigaciones** contra esto, como el canary value. Si el “usuario” ya tomó control del IP, PIE y ASLR dificultan las cosas.

# Alocación dinámica

(To be continued)