

# Princeton University

Computer Science 217: Introduction to Programming Systems



## Data Structures

1

Help you learn (or refresh your memory) about:

- Common data structures: linked lists and hash tables

Why? Deep motivation:

- Common data structures serve as “high level building blocks”
  - A power programmer:
    - Rarely creates programs from scratch
    - Often creates programs using high level building blocks

Why? Shallow motivation:

- Provide background pertinent to Assignment 3
  - ... esp. for those who have not taken COS 226

3

## “Programming in the Large” Steps

### Design & Implement

- Program & programming style (done)
- Common data structures and algorithms <-- we are here
- Modularity
- Building techniques & tools (done)

### Debug

- Debugging techniques & tools (done)

### Test

- Testing techniques (done)

### Maintain

- Performance improvement techniques & tools

2

## Goals of this Lecture



Help you learn (or refresh your memory) about:

- Common data structures: linked lists and hash tables

Why? Deep motivation:

- A power programmer:
  - Rarely creates programs from scratch
  - Often creates programs using high level building blocks

Why? Shallow motivation:

- Provide background pertinent to Assignment 3
  - ... esp. for those who have not taken COS 226

## Common Task



Maintain a collection of key/value pairs

- Each key is a string; each value is an int
- Unknown number of key-value pairs

Examples

- (student name, grade)
  - (“john smith”, 84), (“jane doe”, 93), (“bill clinton”, 81)
  - (baseball player, number)
    - (“Ruth”, 3), (“Gehrig”, 4), (“Mantle”, 7)
    - (variable name, value)
      - (“maxLength”, 2000), (“i”, 7), (“j”, -10)

4

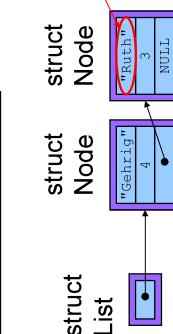
## Agenda

Linked lists

Hash tables

Hash table issues

```
struct Node  
{  
    const char *key;  
    int value;  
    struct Node *next;  
};  
  
struct List  
{  
    struct Node *first;  
};
```



Really this is the address at which “Ruth” resides

5



## Linked List Data Structure



6

## Linked List Data Structure



## Linked List Algorithms



### Create

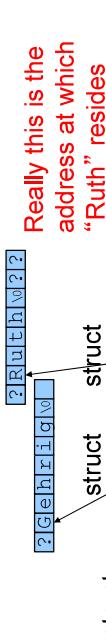
- Allocate List structure; set first to NULL
- Performance: O(1)  $\Rightarrow$  fast

### Add (no check for duplicate key required)

- Insert new node containing key/value pair at front of list
- Performance: O(1)  $\Rightarrow$  fast

### Add (check for duplicate key required)

- Traverse list to check for node with duplicate key
- Insert new node containing key/value pair into list
- Performance: O(n)  $\Rightarrow$  slow



8

## Linked List Algorithms

### Search

- Traverse the list, looking for given key
  - Stop when key found, or reach end
  - Performance: O(n)  $\Rightarrow$  slow
- Free
- Free Node structures while traversing
  - Free List structure
  - Performance: O(n)  $\Rightarrow$  slow
- Would it be better to keep the nodes sorted by key?



## Agenda

### Linked lists

### Hash tables

### Hash table issues

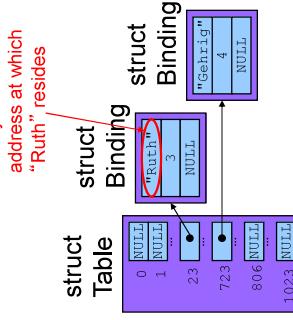
10

## Hash Table Data Structure



### Array of linked lists

```
enum {BUCKET_COUNT = 1024};  
struct Binding  
{  
    const char *key;  
    int value;  
    struct Binding *next;  
};  
struct Table  
{  
    struct Binding *buckets [BUCKET_COUNT];  
};
```

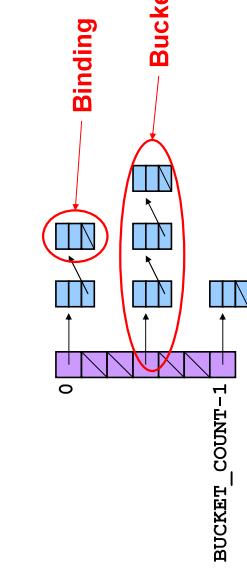


Hash function maps given key to an integer

Mod integer by BUCKET\_COUNT to determine proper bucket

11

## Hash Table Data Structure



12

## Hash Table Example

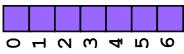
Example: BUCKET\_COUNT = 7  
Add (if not already present) bindings with these keys:  
• the, cat, in, the, hat

13



## Hash Table Example (cont.)

First key: "the"  
• hash("the") = 965156977; 965156977 % 7 = 1  
Search buckets [1] for binding with key "the"; not found

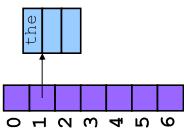


14



## Hash Table Example (cont.)

Add binding with key "the" and its value to buckets [1]

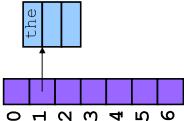


15



## Hash Table Example (cont.)

Second key: "cat"  
• hash("cat") = 3895848756; 3895848756 % 7 = 2  
Search buckets [2] for binding with key "cat"; not found

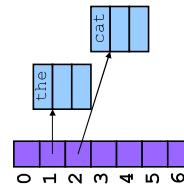


16



## Hash Table Example (cont.)

Third key: "in"  
• hash("in") = 6888005; 6888005 % 7 = 5  
Search buckets [5] for binding with key "in"; not found

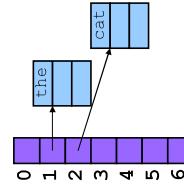


17



## Hash Table Example (cont.)

Add binding with key "cat" and its value to buckets [2]



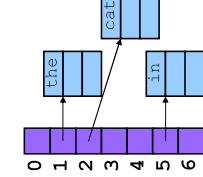
18



## Hash Table Example (cont.)



Add binding with key "in" and its value to buckets [5]



19

## Hash Table Example (cont.)

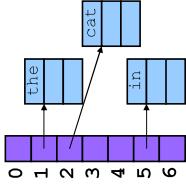


Fourth word: "the"

- $\text{hash}(\text{"the"}) = 965156977; 965156977 \% 7 = 1$

Search buckets [1] for binding with key "the"; found it!

- Don't change hash table



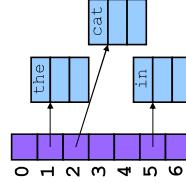
20

## Hash Table Example (cont.)



Add binding with key "hat" and its value to buckets [2]

- At front or back? Doesn't matter
- Inserting at the front is easier, so add at the front



21

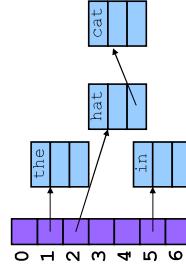
## Hash Table Example (cont.)



Fifth key: "hat"

- $\text{hash}(\text{"hat"}) = 8655559739; 8655559739 \% 7 = 2$

Search buckets [2] for binding with key "hat"; not found



22

## Hash Table Algorithms



### Create

- Allocate Table structure; set each bucket to `NULL`
- Mod by `BUCKET_COUNT` to determine proper bucket

### Add

- Hash the given key
- Mod by `BUCKET_COUNT` to determine proper bucket
- Traverse proper bucket to make sure no duplicate key
- Insert new binding containing key/value pair into proper bucket
- Performance:  $O(1) \Rightarrow \text{fast}$

Is the add performance always fast?

### Search

- Hash the given key
- Mod by `BUCKET_COUNT` to determine proper bucket
- Traverse proper bucket, looking for binding with given key
- Stop when key found, or reach end
- Performance:  $O(1) \Rightarrow \text{fast}$

### Free

- Traverse each bucket, freeing bindings
- Free Table structure
- Performance:  $O(n) \Rightarrow \text{slow}$

Is the search performance always fast?

23

## Hash Table Algorithms



24

## Agenda

- Linked lists
- Hash tables
- Hash table issues



## How Many Buckets?

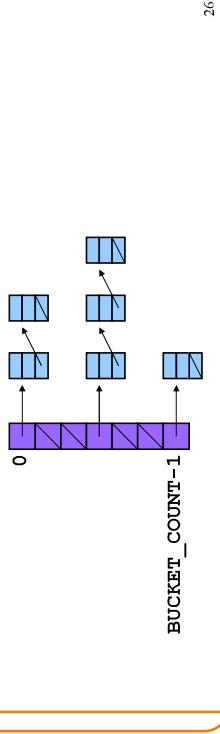
Many!

- Too few  $\Rightarrow$  large buckets  $\Rightarrow$  slow add, slow search

But not too many!

- Too many  $\Rightarrow$  memory is wasted

This is OK:



25

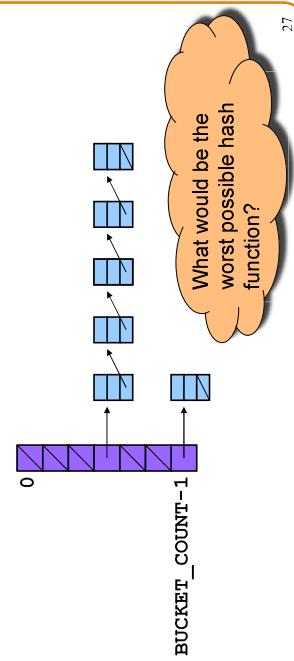
26

## What Hash Function?

Should distribute bindings across the buckets well

- Distribute bindings over the range 0, 1, ..., `BUCKET_COUNT`-1
- Distribute bindings evenly to avoid very long buckets

This is not so good:



27

## How to Hash Strings?

Simple hash schemes don't distribute the keys evenly enough

- Number of characters, mod `BUCKET_COUNT`
- Sum the numeric codes of all characters, mod `BUCKET_COUNT`
- ...

A reasonably good hash function:

- Weighted sum of characters  $s_i$  in the string  $s$ 
  - $(\sum a^i s_i) \text{ mod } BUCKET\_COUNT$
  - Best if  $a$  and `BUCKET_COUNT` are relatively prime
    - E.g.,  $a = 65599, BUCKET\_COUNT = 1024$

Footnote: I originally designed this homework so that `BUCKET_COUNT` is a prime number. In 2016 I wondered, "wouldn't it work just as well if  $a$  and `BUCKET_COUNT` are just relatively prime?" Measurements show no: using a prime number of buckets leads to more even distribution of bucket contents."

28

## How to Hash Strings?

Yielding this function

```
size_t hash(const char *s, size_t bucketCount)
{
    size_t i;
    size_t h = 0;
    for (i=0; s[i]!='\0'; i++)
        h = h * 65599 + (size_t)s[i];
    return h % bucketCount;
}
```

h =  $\Sigma 65599^i * s_i$   
h =  $65599^0 * s_0 + 65599^1 * s_1 + 65599^2 * s_2 + 65599^3 * s_3 + 65599^4 * s_4$   
Direction of traversal of  $s$  doesn't matter, so...  
h =  $65599^0 * s_4 + 65599^1 * s_3 + 65599^2 * s_2 + 65599^3 * s_1 + 65599^4 * s_0$   
h =  $65599^4 * s_0 + 65599^3 * s_1 + 65599^2 * s_2 + 65599^1 * s_3 + 65599^0 * s_4$   
h = (((((s\_0) \* 65599 + s\_1) \* 65599 + s\_2) \* 65599 + s\_3) \* 65599) + s\_4

29

30

## How to Protect Keys?

Suppose `Table_add()` function contains this code:

```
void Table_add(struct Table *t, const char *key, int value)
{
    ...
    struct Binding *p =
        (struct Binding*)malloc(sizeof(struct Binding));
    p->key = key;
    ...
}
```

31

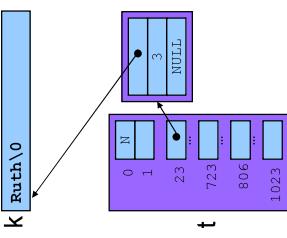


## How to Protect Keys?

Problem: Consider this calling code:

```
struct Table *t;
char k[100] = "Ruth";
...
Table_add(t, k, 3);
```

32



## How to Protect Keys?

Problem: Consider this calling code:

```
struct Table *t;
char k[100] = "Ruth";
...
Table_add(t, k, 3);
strcpy(k, "Gehrig");
```

What happens if the client searches t for "Ruth"? For Gehrig?

33

## How to Protect Keys?

Problem: Consider this calling code:

```
void Table_add(struct Table *t, const char *key, int value)
{
    ...
    struct Binding *p =
        (struct Binding*)malloc(sizeof(struct Binding));
    p->key = (const char*)malloc(strlen(key) + 1);
    strcpy((char*)p->key, key);
    ...
}
```

Why add 1?

34

## How to Protect Keys?

Solution: `Table_add()` saves a **defensive copy** of the given key

```
void Table_add(struct Table *t, const char *key, int value)
{
    ...
    struct Binding *p =
        (struct Binding*)malloc(sizeof(struct Binding));
    p->key = (const char*)malloc(strlen(key) + 1);
    strcpy((char*)p->key, key);
    ...
}
```

## How to Protect Keys?

Now consider same calling code:

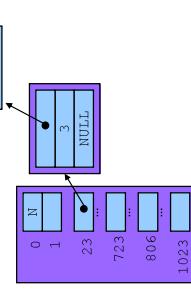
```
struct Table *t;
char k[100] = "Ruth";
...
Table_add(t, k, 3);
strcpy(k, "Gehrig");
```

Hash table is  
not corrupted

## How to Protect Keys?

Now consider same calling code:

```
struct Table *t;
char k[100] = "Ruth";
...
Table_add(t, k, 3);
```

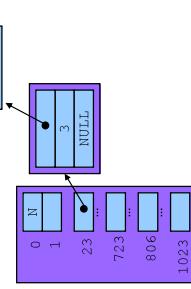


35

## How to Protect Keys?

Now consider same calling code:

```
struct Table *t;
char k[100] = "Ruth";
...
Table_add(t, k, 3);
```



36

## Who Owns the Keys?

Then the **hash table** **owns** its keys

- That is, the hash table owns the memory in which its keys reside
- `Hash_free()` function must free the memory in which the key resides

## Summary

Common data structures and associated algorithms

- Linked list
  - (Maybe) fast add
  - Slow search
- Hash table
  - (Potentially) fast add
  - (Potentially) fast search
- Very common

38

### Hash table issues

- Hashing algorithms
- Defensive copies
- Key ownership



Common data structures and associated algorithms

- Linked list
  - (Maybe) fast add
  - Slow search
- Hash table
  - (Potentially) fast add
  - (Potentially) fast search
- Very common

38

## Princeton University

Computer Science 217: Introduction to Programming Systems



## Debugging (Part 2)



39



## “Programming in the Large” Steps



### Design & Implement

- Program & programming style (**done**)
- Common data structures and algorithms
- Modularity
- Building techniques & tools (**done**)

### Test

- Testing techniques (**done**)

### Debug

- Debugging techniques & tools <-- we are still here

### Maintain

- Performance improvement techniques & tools

40



## “Programming in the Large” Steps



### Design & Implement

- Program & programming style (**done**)
- Common data structures and algorithms
- Modularity
- Building techniques & tools (**done**)

### Test

- Testing techniques (**done**)

### Debug

- Debugging techniques & tools <-- we are still here

### Maintain

- Performance improvement techniques & tools

40



## Goals of this Lecture



Help you learn about:

- Debugging strategies & tools related to dynamic memory management (DMM) \*

Why?

- Many bugs occur in code that does DMM
- DMM errors can be difficult to find
- DMM error in one area can manifest itself in a distant area
- A power programmer knows a wide variety of DMM debugging strategies
- A power programmer knows about tools that facilitate DMM debugging

\* Management of heap memory via `malloc()`, `calloc()`,  
`realloc()`, and `free()`

41



### (9) Look for common DMM bugs

- (10) Diagnose seg faults using `gdb`
- (11) Manually inspect malloc calls
- (12) Hard-code malloc calls
- (13) Comment-out free calls
- (14) Use MemInfo
- (15) Use Valgrind

42

## Look for Common DMM Bugs

Some of our favorites:

```
int *p; /* value of p undefined */
*p = somevalue;

char *p; /* value of p undefined */
fgets(p, 1024, stdin);

int *p;
p = (int*)malloc(sizeof(int));
*p = 5;
free(p);
*p = 6;
```

43

What are  
the  
errors?



## Look for Common DMM Bugs

Some of our favorites:

```
int *p;
...
p = (int*)malloc(sizeof(int));
*p = 5;
...
free(p);
...
free(p);
```

44

What are  
the  
errors?

## Agenda

- (9) Look for common DMM bugs
- (10) Diagnose seg faults using gdb
- (11) Manually inspect malloc calls
- (12) Hard-code malloc calls
- (13) Comment-out free calls
- (14) Use Meminfo
- (15) Use Valgrind

45

## Diagnose Seg Faults Using GDB

Segmentation fault => make it happen in gdb

- Then issue the gdb where command
  - Output will lead you to the line that caused the fault
    - But that line may not be where the error resides!

46



## Agenda

- (9) Look for common DMM bugs
- (10) Diagnose seg faults using gdb
- (11) **Manually inspect malloc calls**
- (12) Hard-code malloc calls
- (13) Comment-out free calls
- (14) Use Meminfo
- (15) Use Valgrind

47

## Manually Inspect Malloc Calls

- Manually inspect each call of malloc()
- Make sure it allocates enough memory

Do the same for calloc() and realloc()

48



# Manually Inspect Malloc Calls

## Agenda



### Some of our favorites:

```
char *s1 = "Hello, world";
char *s2;
s2 = (char*)malloc(strlen(s1));
strcpy(s2, s1);

char *s1 = "Hello";
char *s2;
s2 = (char*)malloc(sizeof(s1));
strcpy(s2, s1);

long double *p;
p = (long double*)malloc(sizeof(long double*));

long double *p;
p = (long double*)malloc(sizeof(long double*));
```

49



## Hard-Code Malloc Calls



Temporarily change each call of `malloc()` to request a large number of bytes

- Say, 10000 bytes
- If the error disappears, then at least one of your calls is requesting too few bytes

Then incrementally restore each call of `malloc()` to its previous form

- When the error reappears, you might have found the culprit

Do the same for `calloc()` and `realloc()`

51

## Comment-Out Free Calls



Temporarily comment-out every call of `free()`

- If the error disappears, then program is
  - Freeing memory too soon, or
  - Freeing memory that already has been freed, or
  - Freeing memory that should not be freed,
  - Etc.

Then incrementally “comment-in” each call of `free()`

- When the error reappears, you might have found the culprit

53

### (9) Look for common DMM bugs

- (10) Diagnose seg faults using gdb
- (11) Manually inspect malloc calls
- (12) Hard-code malloc calls**
- (13) Comment-out free calls
- (14) Use Meminfo
- (15) Use Valgrind

50

## Agenda



### (9) Look for common DMM bugs

- (10) Diagnose seg faults using gdb
- (11) Manually inspect malloc calls
- (12) Hard-code malloc calls**
- (13) Comment-out free calls**
- (14) Use Meminfo
- (15) Use Valgrind

52

## Agenda



- (9) Look for common DMM bugs
- (10) Diagnose seg faults using gdb
- (11) Manually inspect malloc calls
- (12) Hard-code malloc calls**
- (13) Comment-out free calls**
- (14) Use Meminfo**
- (15) Use Valgrind

54

# Use Meminfo

## Use the **Meminfo** tool

- Simple tool
  - Initial version written by Dondero
  - Current version written by COS 217 alumnus RJ Liljestrom
  - Reports errors after program execution
    - Memory leaks
    - Some memory corruption
    - User-friendly output

Appendix 1 provides example buggy programs

Appendix 2 provides Meminfo analyses

55



# Agenda

- (9) Look for common DMM bugs
- (10) Diagnose seg faults using gdb
- (11) Manually inspect malloc calls
- (12) Hard-code malloc calls
- (13) Comment-out free calls
- (14) Use Meminfo
- (15) Use Valgrind**

56

# Use Valgrind

## Use the **Valgrind** tool

- Complex tool
  - Written by multiple developers, worldwide
  - See [www.valgrind.org](http://www.valgrind.org)
  - Reports errors during program execution
    - Memory leaks
    - Multiple frees
    - Dereferences of dangling pointers
    - Memory corruption
    - Comprehensive output
    - But not always user-friendly

57



# Use Valgrind

Appendix 1 provides example buggy programs

Appendix 3 provides Valgrind analyses

58



# Use Valgrind

- Complex tool
  - Written by multiple developers, worldwide
  - See [www.valgrind.org](http://www.valgrind.org)
  - Reports errors during program execution
    - Memory leaks
    - Multiple frees
    - Dereferences of dangling pointers
    - Memory corruption
    - Comprehensive output
    - But not always user-friendly

59

# Summary

Strategies and tools for debugging the DMM aspects of your code:

- Look for common DMM bugs
- Diagnose seg faults using gdb
- Manually inspect malloc calls
- Hard-code malloc calls
- Comment-out free calls
- Use Meminfo
- Use Valgrind



# Appendix 1: Buggy Programs

leak.c

```
1. #include <stdio.h>
2. #include <stdlib.h>
3. int main(void)
4. {
5.     int *pi;
6.     *pi = 5;
7.     printf("%d\n", *pi);
8.     pi = (int*)malloc(sizeof(int));
9.     *pi = 6;
10.    printf("%d\n", *pi);
11.    free(pi);
12.    return 0;
13. }
```

Memory leak:

Memory allocated at line 5 is leaked

59

60



## Appendix 3: Valgrind

### Valgrind can detect memory leaks:

```
$ gcc217 leak.c -o leak
$ valgrind leak
==31921== Memcheck, a memory error detector
==31921== Copyright (C) 2002-2012, and GNU GPL'd, by Julian Seward et al.
==31921== Using Valgrind-3.8.1 and LibEX: rerun with -h for copyright info
==31921== Command: leak
==31921==

5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
```

## Appendix 3: Valgrind

### Valgrind can detect memory leaks:

```
$ valgrind --leak-check=full leak
==476== Memcheck, a memory error detector
==476== Copyright (C) 2002-2012, and GNU GPL'd, by Julian Seward et al.
==476== Using Valgrind-3.8.1 and LibEX: rerun with -h for copyright info
==476== Command: leak
==476==

5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
```

## Appendix 3: Valgrind

### Valgrind can detect multiple frees:

```
$ gcc217 doublefree.c -o doublefree
$ valgrind doublefree
==31951== Memcheck, a memory error detector
==31951== Copyright (C) 2002-2012, and GNU GPL'd, by Julian Seward et al.
==31951== Using Valgrind-3.8.1 and LibEX: rerun with -h for copyright info
==31951== Command: doublefree
==31951== Invalid free() / delete[] / realloc()
==31951== at 0x4006380: free (vg_replace_malloc.c:446)
==31951== by 0x4002A5: main (doublefree.c:9)
==31951== Address 0x4c2a000 is 0 bytes inside a block of size 4 freed
==31951== at 0x4006380: free (vg_replace_malloc.c:446)
==31951== by 0x400599: main (doublefree.c:8)
==31951==

5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
```

## Appendix 3: Valgrind

### Valgrind can detect dereferences of dangling pointers:

```
$ gcc217 danglingptr.c -o danglingptr
$ valgrind --danglingptr=danglingptr
==316== Memcheck, a memory error detector
==316== Copyright (C) 2002-2012, and GNU GPL'd, by Julian Seward et al.
==316== Using Valgrind-3.8.1 and LibEX: rerun with -h for copyright info
==316== Command: danglingptr
==316== Invalid read of size 4
==316== at 0x400599: main (danglingptr.c:9)
==316== Address 0x4ac0d40 is 0 bytes inside a block of size 4 freed
==316== at 0x4006380: free (vg_replace_malloc.c:446)
==316== by 0x400599: main (danglingptr.c:8)
==316==

5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
```

## Appendix 3: Valgrind

### Valgrind can detect memory corruption (cont.):

```
$ valgrind --error-suppress=danglingptr --leak-check=full --show-leak-kinds=all --track-origins=yes --log-file=valgrind.log
==436== Invalid write of size 4
==436== at 0x40063E: main (toosmall.c:6)
==436== Address 0x4c2a000 is 0 bytes inside a block of size 1 alloc'd
==436== at 0x400592E: malloc (vg_replace_malloc.c:270)
==436== by 0x400565: main (toosmall.c:5)
==436== Invalid read of size 4
==436== at 0x400578: main (toosmall.c:7)
==436== Address 0x4c2a000 is 0 bytes inside a block of size 1 alloc'd
==436== at 0x400592E: malloc (vg_replace_malloc.c:270)
==436== by 0x400565: main (toosmall.c:5)
==436==

5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
```

## Appendix 3: Valgrind

### Continued from previous slide

```
$ valgrind --error-suppress=danglingptr --leak-check=full --show-leak-kinds=all --track-origins=yes --log-file=valgrind.log
==436== Invalid write of size 4
==436== at 0x40063E: main (toosmall.c:6)
==436== Address 0x4c2a000 is 0 bytes inside a block of size 1 alloc'd
==436== at 0x400592E: malloc (vg_replace_malloc.c:270)
==436== by 0x400565: main (toosmall.c:5)
==436== Invalid read of size 4
==436== at 0x400578: main (toosmall.c:7)
==436== Address 0x4c2a000 is 0 bytes inside a block of size 1 alloc'd
==436== at 0x400592E: malloc (vg_replace_malloc.c:270)
==436== by 0x400565: main (toosmall.c:5)
==436== All heap blocks were freed -- no leaks are possible
==336== For counts of detected and suppressed errors, rerun with: --v
==336== ERROR SUMMARY: 1 errors from 1 contexts (suppressed: 6 from 6)
==336==
```

## Appendix 3: Valgrind

### Continued on next slide

## Appendix 3: Valgrind



### Valgrind caveats:

- Not intended for programmers who are new to C
  - Messages may be cryptic
- Suggestion:
  - Observe line numbers referenced by messages
  - Study code at those lines
  - Infer meanings of messages