

N2624: Introduction for discussion of N2577 Provenance-aware Memory Object Model for C Draft Technical Specification

adapted from N2378

Jens Gustedt¹ Peter Sewell² Kayvan Memarian²
Victor B. F. Gomes^{2*} Martin Uecker³

¹INRIA and ICube, Université de Strasbourg, France

²University of Cambridge, UK *when this work was done

³University Medical Center, Göttingen, Germany

Context: previous discussions (selected)

WG14 online meeting, 2020-11

N2577: Provenance-aware Memory Object Model for C, Draft Technical Specification: Jens Gustedt, Peter Sewell, Kayvan Memarian, Victor B. F. Gomes, Martin Uecker

WG21 Cologne meeting, 2019-07

C/C++ Memory Object Model Papers - Introduction (P1797R0): Peter Sewell

Effective Types: Examples (P1796R0): Peter Sewell, Kayvan Memarian, Victor B. F. Gomes, Jens Gustedt, Hubert Tong

WG14 London meeting, 2019-04

N2378: C provenance semantics: slides (extracts from N2363): Peter Sewell, Kayvan Memarian, Victor B. F. Gomes, Jens Gustedt, and Martin Uecker

N2362: Moving to a provenance-aware memory object model for C: Jens Gustedt, Peter Sewell, Kayvan Memarian, Victor B. F. Gomes, and Martin Uecker

N2363: C provenance semantics: examples: Peter Sewell, Kayvan Memarian, Victor B. F. Gomes, Jens Gustedt, and Martin Uecker

N2364: C provenance semantics: detailed semantics (for PNVI-plain, PNVI address-exposed, PNVI address-exposed user-disambiguation, and PVI models): Peter Sewell, Kayvan Memarian, and Victor B. F. Gomes

N2369: Pointer lifetime-end zap: Paul E. McKenney, Maged Michael, and Peter Sewell

WG14 Pittsburgh meeting, 2018-10

n2294: C Memory Object Model Study Group: Progress ReportPeter Sewell

n2263: Clarifying Pointer Provenance v4Kayvan Memarian, Victor Gomes, Peter Sewell

WG14 Brno meeting, 2018-04 (CMOM SG created)

n2223: Clarifying the C Memory Object Model: Introduction to N2219 - N2222Kayvan Memarian, Victor Gomes, Peter Sewell superseded by v4

n2219: Clarifying Pointer Provenance (Q1-Q20) v3Kayvan Memarian, Victor Gomes, Peter Sewell superseded by v4

n2220: Clarifying Trap Representations (Q47) v3Kayvan Memarian, Victor Gomes, Peter Sewell superseded by the above v4

n2221: Clarifying Unspecified Values (Q48-Q59) v3Kayvan Memarian, Victor Gomes, Peter Sewell superseded by v4

n2222: Further Pointer Issues (Q21-Q46)Kayvan Memarian, Victor Gomes, Peter Sewell superseded by v4

WG14 Pittsburgh meeting, 2016-10

n2089: Clarifying Unspecified Values (Draft Defect Report or Proposal for C2x)Kayvan Memarian and Peter Sewell

n2090: Clarifying Pointer Provenance (Draft Defect Report or Proposal for C2x)Kayvan Memarian and Peter Sewell

n2091: Clarifying Trap Representations (Draft Defect Report or Proposal for C2x)Kayvan Memarian and Peter Sewell

WG14 London meeting, 2016-04

n2012: Clarifying the C memory object modelKayvan Memarian and Peter Sewell

n2013: C Memory Object and Value Semantics: The Space of de facto and ISO StandardsDavid Chisnall, Justus Matthesen, Kayvan Memarian, Peter Sewell, Robert N. M. Watson

n2014: What is C in Practice? (Cerberus Survey v2): Analysis of ResponseKayvan Memarian and Peter Sewell

n2015: What is C in practice? (Cerberus survey v2): Analysis of Responses - with CommentsKayvan Memarian and Peter Sewell

Academic papers

Exploring C Semantics and Pointer Provenance (in POPL 2019, and as n2311)

Into the depths of C: elaborating the de facto standards (in PLDI 2016)

Elsewhere and Previously

WG21 p0137r1: Core Issue 1776: Replacement of class objects containing reference members (in C++17) (2016-06)

WG21 p0593r3: Implicit creation of objects for low-level object manipulation (2019-01)

In OOPSLA 2018: Reconciling High-level Optimizations and Low-level Code in LLVM (2018-11)

n1818 / DR451: Defect Report 451 (2014-04)

n1637: Subtleties of the ANSI/ISO C standard (2012-09)

DR260: indeterminate values and identical representations (2004-09)

mail 9350: What is an Object in C Terms? (2001-09)

This meeting (WG14 online, 2020-11)

Previously established a **well-developed proposal for pointer provenance**:, presented in WG14 London 2019-04, WG21 Pittsburgh 2018-10, GCC Cauldron 2018, and EuroLLVM 2018

Here, following earlier WG14 discussion: we have a draft Technical Specification

N2577: Provenance-aware Memory Object Model for C, Draft Technical Specification. Jens Gustedt, Peter Sewell, Kayvan Memarian, Victor B. F. Gomes, Martin Uecker

Incorporates the content of

- ▶ **Examples:**

N2363: C provenance semantics: examples. Peter Sewell, Kayvan Memarian, Victor B. F. Gomes, Jens Gustedt, and Martin Uecker [these slides are mostly extracts from this]

- ▶ **Mathematical semantics:**

N2364: C provenance semantics: detailed semantics (for PNVI-plain, PNVI address-exposed, PNVI address-exposed user-disambiguation, and PVI models). Peter Sewell, Kayvan Memarian, and Victor B. F. Gomes

- ▶ **Proposed standard text diff:**

N2362: Moving to a provenance-aware memory object model for C. Jens Gustedt, Peter Sewell, Kayvan Memarian, Victor B. F. Gomes, and Martin Uecker

First two exactly as before; the last updated as a change-highlighted diff w.r.t. ISO/IEC 9899:2018

Also:

- ▶ **Executable Web-GUI semantics in Cerberus:**

<http://cerberus.cl.cam.ac.uk/cerberus>

Scope

The draft TS only addresses basic pointer provenance.

It does not address many other open memory-object-model questions:

- ▶ uninitialised values
- ▶ padding
- ▶ effective types
- ▶ subobject provenance
- ▶ pointer lifetime-end zap
- ▶ ...various

All those are important, and we've considered them in detail in other papers (not all yet with clear solutions). But to work within WG14 bandwidth limits, we've focussed on the basic provenance for now.

Open question: will we stick with just that for the TS, or incorporate more?

Depends how quickly we can process this.

Upcoming planned memory object model study group work on uninitialised values and padding.

Process

New Work Item Proposal ballot closed yesterday and passed (needed five countries; got seven) – thanks to David K. and all for navigating that.

Now WG14 has two years, until December 1, 2022, to submit a final draft for ballot (DTS, Draft Technical Specification). Ballot can be earlier if we choose.

The publication deadline is one year after that, which gives us time to navigate the ISO publication process.

Basic pointer provenance

Recall: pointers are typically simple concrete addresses at runtime, but compilers do *provenance-based alias analysis*:

```
// provenance_basic_global_yx.c
1 #include <stdio.h>
2 #include <string.h>
3 int y=2, x=1;
4 int main() {
5     int *p = &x;
6     int *q = &y;
7     p=p+1;
8     printf("Addresses: p=%p q=%p\n", (void*)p, (void*)q);
9     if (memcmp(&p, &q, sizeof(p)) == 0) {
10         *p = 11; // does this have undefined behaviour?
11         printf("x=%d y=%d *p=%d *q=%d\n", x, y, *p, *q);
12     }
13 }
```

Clang 6.0 -O2	x=1 y=11 *p=11 *q=11
GCC 8.1 -O2	x=1 y=2 *p=11 *q=2
ICC 19 -O2	x=1 y=2 *p=11 *q=11 (with x and y swapped)

Basic pointer provenance

To make that GCC and ICC compiler optimisation legal w.r.t. the standard, this program has to be deemed to have UB, so we have to recognise that the `p` in `*p=11` is a one-past pointer, even though it has the same numeric address as `&y`.

DR260 CR (2001) hints at this:

“Implementations are permitted to track the origins of a bit-pattern [...]. They may also treat pointers based on different origins as distinct even though they are bitwise identical.”

but it was never incorporated in the standard text, and it gives no more detail. That leaves unclear whether some programming idioms are allowed or not, and what compiler alias analysis and optimisation are allowed to do.

Our proposal

Our proposal clarifies this. It reconciles existing C programming practice, compiler implementation practice, and the standard text, as best we can, with a well-defined and reasonably simple semantics.

We aim to be conservative with respect to all those – as far as possible, the proposal is capturing the status quo in the specification. The proposal doesn't involve any new features or change to the language syntax.

Our proposal: the basic idea

We associate a *provenance* $@i$ with every pointer value in the abstract machine, identifying the original storage instance it's derived from (if any, or $@\text{empty}$ otherwise).

- ▶ On every allocation (for static, thread, automatic, and allocated storage durations), the abstract machine chooses a fresh storage instance ID i (unique across the entire execution), and the resulting pointer value carries that as its provenance $@i$.
- ▶ Provenance is preserved by pointer arithmetic that adds or subtracts an integer to a pointer.
- ▶ At any access via a pointer value, its numeric address must be consistent with its provenance, with undefined behaviour otherwise

That UB is what licences compilers to do provenance-aware alias analysis.

Our proposal: the basic idea

The screenshot shows a debugger window with two panes. The left pane displays the source code for `provenance_basic_global_yx.c`. The right pane shows a memory diagram with two pointers, `p` and `q`, and their respective target memory locations.

```
1 #include <stdio.h>
2 #include <string.h>
3 int y=2, x=1;
4 int main() {
5     int *p = &x;
6     int *q = &y;
7     p=p+1;
8     if (memcmp(&p, &q, sizeof(p)) == 0) {
9         *p = 11; // does this have UB?
10        printf("x=%d y=%d *p=%d
11               *q=%d\n", x, y, *p, *q);
12    }
13 }
```

The memory diagram shows:

- Pointer `p`: `p: signed int* [@4, 0xfffffd0]`. Its value is `@2, 0xfffffe0`. An arrow points to a memory location `x: signed int [@2, 0xfffffe0]` containing the value `1`.
- Pointer `q`: `q: signed int* [@5, 0xfffffc8]`. Its value is `@1, 0xfffffe4`. An arrow points to a memory location `y: signed int [@1, 0xfffffe4]` containing the value `2`.

Note the storage instance IDs `@i` of the allocations and as part of the pointer values.

(try it live at https://cerberus.cl.cam.ac.uk/cerberus?defacto/provenance_basic_global_yx.c)

Our proposal: the basic idea

The screenshot shows a debugger window with the following components:

- Code Editor:** Shows the C source code for `provenance_basic_global_yx.c`. The code includes `<stdio.h>` and `<string.h>`, defines `int y=2, x=1;`, and a `main` function. In `main`, `int *p = &x;` and `int *q = &y;` are declared. `p` is incremented (`p=p+1;`). A conditional block checks `memcmp(&p, &q, sizeof(p)) == 0`. Inside, `*p = 11;` is executed, and `printf` prints `x=%d y=%d *p=%d *q=%d`. The code ends with `}`.
- Memory View:** Shows the state of memory:
 - `x: signed int [@2, 0xfffffe0]`: Contains the value `1`.
 - `p: signed int* [@4, 0xfffffd0]`: Contains the address `@2, 0xfffffe4`. A red arrow points from this box to the `y` box.
 - `q: signed int* [@5, 0xfffffc8]`: Contains the address `@1, 0xfffffe4`. A black arrow points from this box to the `y` box.
 - `y: signed int [@1, 0xfffffe4]`: Contains the value `2`.

after the `p=p+1`, `p` has the address of `y` (in this execution), but it still has the provenance (`@2`) of `x`

Our proposal: the basic idea

The screenshot shows a debugger window for the file `provenance_basic_global_yx.c`. The code in the left pane is as follows:

```
1 #include <stdio.h>
2 #include <string.h>
3 int y=2, x=1;
4 int main() {
5     int *p = &x;
6     int *q = &y;
7     p=p+1;
8     if (memcmp(&p, &q, sizeof(p)) == 0) {
9         *p = 11; // does this have UB?
10        printf("x=%d y=%d *p=%d
11               *q=%d\n", x, y, *p, *q);
12    }
13 }
```

The console output shows:

```
1 Unsuccessful termination of this executi
2 Undefined: [out of bounds pointer at m
```

The memory dump on the right shows the following state:

- `x: signed int` at address `@2 exp, 0xfffffe0` contains the value `1`.
- `p: signed int*` at address `@4, 0xfffffd0` points to `@2, 0xfffffe4`.
- `q: signed int*` at address `@5, 0xfffffc8` points to `@1, 0xfffffe4`.
- `y: signed int` at address `@1 exp, 0xfffffe4` contains the value `2`.

Red and black arrows indicate that both `*p` and `*q` point to the memory location `@1, 0xfffffe4`, which is the location of `y`.

at the `*p=11` access, the address is not within the footprint of the allocation with that provenance, so the access is UB, as required

So far so good, but...

C provides many other ways to construct pointer values:

- ▶ casts of pointers to integer types and back, possibly with integer arithmetic
- ▶ copying pointer values with `memcpy`
- ▶ manipulation of the representation bytes of pointers, e.g. via `char*` accesses
- ▶ type punning between pointer and integer values
- ▶ I/O, using either `fprintf/fscanf` and the `%p` format, `fwrite/fread` on the pointer representation bytes, or pointer/integer casts and integer I/O
- ▶ copying pointer values with `realloc`
- ▶ constructing pointer values that embody knowledge established from linking, and from constants that represent the addresses of memory-mapped devices.

We have to address all these, and the impact on optimisation.

Design options

- ▶ **PVI**: track provenance via integer computation ([n2090](#), [n2263](#))
Complex, poor algebraic properties, not good fit with implementation

Design options

- ▶ **PVI**: track provenance via integer computation ([n2090](#), [n2263](#))
Complex, poor algebraic properties, not good fit with implementation
- ▶ **PNVI-plain**: don't track provenance via integers. Instead, at integer-to-pointer cast points, check whether the given address points within a live object and, if so, recreate the corresponding provenance.

Design options

- ▶ **PVI**: track provenance via integer computation ([n2090](#), [n2263](#))
Complex, poor algebraic properties, not good fit with implementation
- ▶ **PNVI-plain**: don't track provenance via integers. Instead, at integer-to-pointer cast points, check whether the given address points within a live object and, if so, recreate the corresponding provenance.
- ▶ **PNVI-exposed-address (PNVI-ae)**: allow integer-to-pointer casts to recreate provenance only for storage instances that have previously been *exposed*, by a cast of a pointer to it to an integer type, by a read (at non-pointer type) of the representation of such a pointer, or by an output of such a pointer using %p.

Design options

- ▶ **PVI**: track provenance via integer computation ([n2090](#), [n2263](#))
Complex, poor algebraic properties, not good fit with implementation
- ▶ **PNVI-plain**: don't track provenance via integers. Instead, at integer-to-pointer cast points, check whether the given address points within a live object and, if so, recreate the corresponding provenance.
- ▶ **PNVI-exposed-address (PNVI-ae)**: allow integer-to-pointer casts to recreate provenance only for storage instances that have previously been *exposed*, by a cast of a pointer to it to an integer type, by a read (at non-pointer type) of the representation of such a pointer, or by an output of such a pointer using %p.
- ▶ **PNVI exposed-address user-disambiguation (PNVI-ae-udi)**: a further refinement to support roundtrip casts, pointer to integer and back, of one-past pointers.
Clear preference for this option from WG14 and WG21 UB group in previous meetings

Next

- ▶ **Idioms:** checking that various desirable idioms work
- ▶ **Implications for optimisation:** checking that various cases are UB, e.g. that function arguments can't alias its local variables
- ▶ **PNVI-plain vs PNVI-ae-*:** is the “exposed” machinery needed?
- ▶ **PNVI-ae vs PNVI-ae-udi:** what about one-past integers?
- ▶ **Experimental checks:** running the examples in an executable model and in GCC/Clang/ICC
- ▶ **The proposed text diff (Jens)**
- ▶ **Precise semantics**

Idioms

Pointer/integer casts

The screenshot shows a debugger window with two panes. The left pane displays C code for a program named 'provenance_roundtrip_via_intptr_t.c'. The code includes `<stdio.h>` and `<inttypes.h>`, declares `int x=1;`, and defines a `main()` function. Inside `main()`, it declares `int *p = &x;`, casts `intptr_t i = (intptr_t)p;`, and declares `int *q = (int *)i;`. A comment `// is this free of UB?` is placed above `*q = 11;`. The `printf` statement prints the values of `*p` and `*q`. The right pane shows a memory dump for the 'Memory' window. It displays three memory locations: `q: signed int* [@5, 0xfffffd0]` (unspecified), `i: intptr_t [@4, 0xfffffd8]` (unspecified), and `p: signed int* [@3, 0xfffffe0]`. The `p` location contains the address `@1, 0xfffffec`, which points to the memory location `x: signed int [@1, 0xfffffec]` containing the value `1`.

```
1 #include <stdio.h>
2 #include <inttypes.h>
3 int x=1;
4 int main() {
5     int *p = &x;
6     intptr_t i = (intptr_t)p;
7     int *q = (int *)i;
8     *q = 11; // is this free of UB?
9     printf("*p=%d *q=%d\n", *p, *q);
10 }
11
```

Memory dump:

- `q: signed int* [@5, 0xfffffd0]` (unspecified)
- `i: intptr_t [@4, 0xfffffd8]` (unspecified)
- `p: signed int* [@3, 0xfffffe0]` (contains `@1, 0xfffffec`)
- `x: signed int [@1, 0xfffffec]` (contains `1`)

This is a simple pointer-to-integer-to-pointer roundtrip; the result should be usable for access.

Pointer/integer casts

The screenshot shows a debugger window with the following C code on the left and a memory dump on the right.

```
1 #include <stdio.h>
2 #include <inttypes.h>
3 int x=1;
4 int main() {
5     int *p = &x;
6     intptr_t i = (intptr_t)p;
7     int *q = (int *)i;
8     *q = 11; // is this free of UB?
9     printf("*p=%d *q=%d\n", *p, *q);
10 }
11
```

The memory dump on the right shows the following state:

- q: signed int* [@5, 0xfffffd0]**: unspecified
- p: signed int* [@3, 0xfffffe0]**: @1, 0xfffffec
- i: intptr_t [@4, 0xfffffd8]**: 0xfffffec
- x: signed int [@1 exp, 0xfffffec]**: 1

The **x** allocation is highlighted in orange and labeled as "exposed". A solid arrow points from the **p** memory location to the **x** allocation, and a dashed arrow points from the **i** memory location to the **x** allocation.

After the pointer-to-integer cast `(intptr_t)p`, the `x` allocation is marked as exposed.

Pointer/integer casts

```
Cerberus ▾  provenance_roundtrip_via_intptr_t.c ▾  File ▾  Model ▾  Views ▾  Step: 12  Forward  Back  Restart  Search ▾
```

```
provenance_roundtrip_via_intptr_t.c
```

```
1 #include <stdio.h>
2 #include <inttypes.h>
3 int x=1;
4 int main() {
5     int *p = &x;
6     intptr_t i = (intptr_t)p;
7     int *q = (int *)i;
8     *q = 11; // is this free of UB?
9     printf("p=%d q=%d\n", *p, *q);
10 }
11
```

Memory x

The diagram illustrates the memory state during the execution of the provided C code. It shows three memory locations:

- p:** signed int* [@3, 0xfffffe0]
- i:** intptr_t [@4, 0xffffd8] containing the value 0xfffffec
- q:** signed int* [@5, 0xffffd0]

A dashed arrow points from the value 0xfffffec in memory location **i** to a highlighted memory location **x:** signed int [@1 exp, 0xfffffec] containing the value 1. This indicates that the integer value stored in **i** is the provenance of the memory location **x**.

so the integer-to-pointer cast `(int*)i`, of an integer within the footprint of `x`, will recover the provenance (`@1`) of `x`

Pointer/integer casts

```
Cerberus ▾ provenance_roundtrip_via_intptr_t.c ▾ File ▾ Model ▾ Views ▾ Step: 15 Forward Back Restart Search ▾
```

```
provenance_roundtrip_via_intptr_t.c
```

```
1 #include <stdio.h>
2 #include <inttypes.h>
3 int x=1;
4 int main() {
5     int *p = &x;
6     intptr_t i = (intptr_t)p;
7     int *q = (int *)i;
8     *q = 11; // is this free of UB?
9     printf("p=%d q=%d\n", *p, *q);
10 }
11
```

Memory

p: signed int* [@3, 0xfffffe0]
@1, 0xfffffec

i: intptr_t [@4, 0xfffffd8]
0xfffffec

q: signed int* [@5, 0xfffffd0]
@1, 0xfffffec

x: signed int [@1 exp, 0xfffffec]
11

and the access `*q=11` is defined behaviour.

Pointer provenance for pointer bit manipulations

Common in practice. For example, assuming `int` has alignment at least 4, the low-order pointer bits are unused, and the implementation-defined pointer/integer conversions are as expected:

```
// provenance_tag_bits_via_uintptr_t_1.c
1 #include <stdio.h>
2 #include <stdint.h>
3 int x=1;
4 int main() {
5     int *p = &x;
6     // cast &x to an integer
7     uintptr_t i = (uintptr_t) p;
8     // set low-order bit
9     i = i | 1u;
10    // cast back to a pointer
11    int *q = (int *) i; // does this have UB?
12    // cast to integer and mask out low-order bits
13    uintptr_t j = ((uintptr_t)q) & ~((uintptr_t)3u);
14    // cast back to a pointer
15    int *r = (int *) j;
16    // are r and p now equivalent?
17    *r = 11; // does this have UB?
18    _Bool b = (r==p); // is this true?
19    printf("x=%i *r=%i (r==p)=%s\n",x,*r,b?"t":"f");
20 }
```

As before, `(uintptr_t)x` will expose `x`, so the `(int*)j` cast will recover the correct provenance, making the access `*r=11` legal.

Inter-object integer arithmetic

Can one move between objects with pointer arithmetic? No.

Can one move between objects with integer arithmetic? Debatable whether this must be supported – we get conflicting reports as to how important it is in practice, e.g. for XOR linked lists.

PNVI-* naturally allows it (if the implementation-defined pointer/integer conversions do).

Inter-object integer arithmetic

```
// pointer_offset_from_int_subtraction_global_yx.c
1 #include <stdio.h>
2 #include <string.h>
3 #include <stdint.h>
4 #include <inttypes.h>
5 int y=2, x=1;
6 int main() {
7     uintptr_t ux = (uintptr_t)&x;
8     uintptr_t uy = (uintptr_t)&y;
9     uintptr_t offset = uy - ux;
10    printf("Addresses: &x=%"PRIuPTR" &y=%"PRIuPTR"\n",
11           " offset=%"PRIuPTR" \n", (unsigned long)ux, (unsigned long)uy, (unsigned long)offset);
12    int *p = (int *) (ux + offset);
13    int *q = &y;
14    if (memcmp(&p, &q, sizeof(p)) == 0) {
15        *p = 11; // is this free of UB?
16        printf("x=%d y=%d *p=%d *q=%d\n", x, y, *p, *q);
17    }
18 }
```

As before: the cast `(uintptr_t)&y` marks `y` as exposed, so the cast `p=(int*)(ux+offset)` can recover the provenance of `y` and make the access `*p=11` legal.

Copying pointer values bitwise, with user-memcpy

C supports manipulation of object representations, e.g. as in the following naive user implementation of memcpy:

```
// pointer_copy_user_dataflow_direct_bytewise.c
1 #include <stdio.h>
2 #include <string.h>
3 int x=1;
4 void user_memcpy(unsigned char* dest,
5                 unsigned char *src, size_t n) {
6     while (n > 0) {
7         *dest = *src;
8         src += 1; dest += 1; n -= 1;
9     }
10 }
11 int main() {
12     int *p = &x;
13     int *q;
14     user_memcpy((unsigned char*)&q,
15               (unsigned char*)&p, sizeof(int *));
16     *q = 11; // is this free of undefined behaviour?
17     printf("*p=%d *q=%d\n", *p, *q);
18 }
```

which constructs a pointer value from copied bytes. This too should be allowed.

Copying pointer values bitwise, with user-memcpy

q: signed int* [@4, 0xfffffd8]

unspecified

p: signed int* [@3, 0xfffffe0]

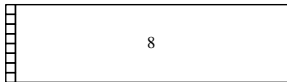
@1, 0xfffffec

x: signed int [@1, 0xfffffec]

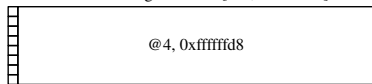
1

Copying pointer values bitwise, with user-memcpy

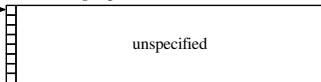
n: size_t [@7, 0xffffffc0]



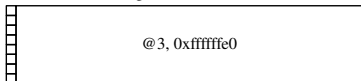
dest: unsigned char* [@5, 0xfffffd0]



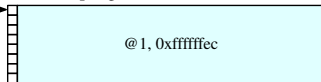
q: signed int* [@4, 0xfffffd8]



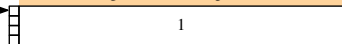
src: unsigned char* [@6, 0xfffffc8]



p: signed int* [@3, 0xfffffe0]

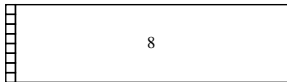


x: signed int [@1 exp, 0xfffffec]

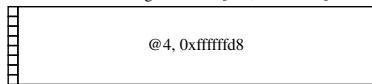


Copying pointer values bitwise, with user-memcpy

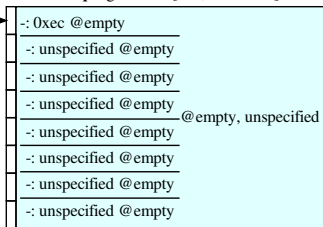
n: size_t [@7, 0xffffffc0]



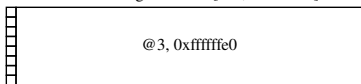
dest: unsigned char* [@5, 0xfffffd0]



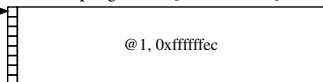
q: signed int* [@4, 0xfffffd8]



src: unsigned char* [@6, 0xfffffc8]



p: signed int* [@3, 0xfffffe0]



x: signed int [@1 exp, 0xfffffec]

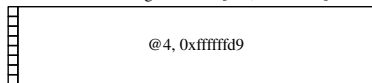


Copying pointer values bitwise, with user-memcpy

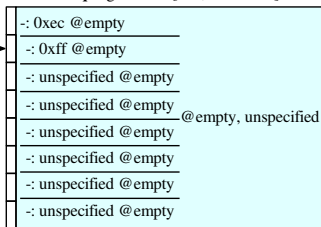
n: size_t [@7, 0xfffffc0]



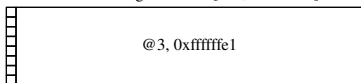
dest: unsigned char* [@5, 0xfffffd0]



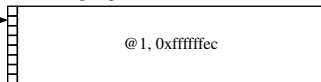
q: signed int* [@4, 0xfffffd8]



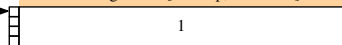
src: unsigned char* [@6, 0xfffffc8]



p: signed int* [@3, 0xfffffe0]

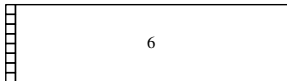


x: signed int [@1 exp, 0xfffffec]

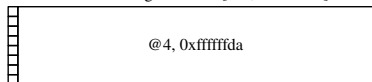


Copying pointer values bitwise, with user-memcpy

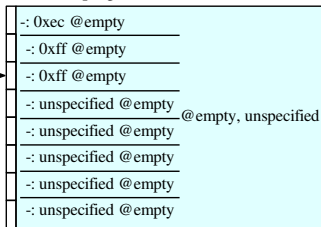
n: size_t [@7, 0xfffffc0]



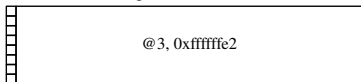
dest: unsigned char* [@5, 0xfffffd0]



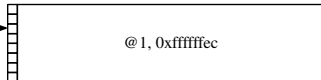
q: signed int* [@4, 0xfffffd8]



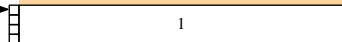
src: unsigned char* [@6, 0xfffffc8]



p: signed int* [@3, 0xfffffe0]



x: signed int [@1 exp, 0xfffffec]

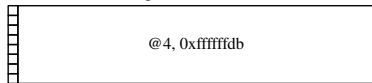


Copying pointer values bitwise, with user-memcpy

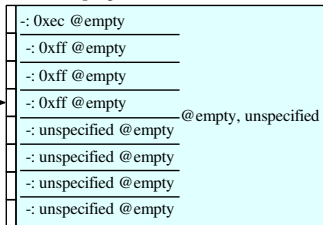
n: size_t [@7, 0xfffffc0]



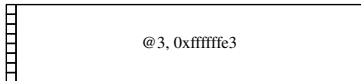
dest: unsigned char* [@5, 0xfffffd0]



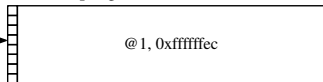
q: signed int* [@4, 0xfffffd8]



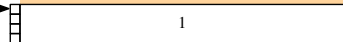
src: unsigned char* [@6, 0xfffffc8]



p: signed int* [@3, 0xfffffe0]

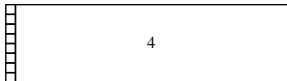


x: signed int [@1 exp, 0xfffffec]

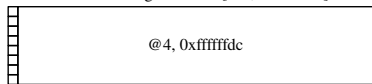


Copying pointer values bitwise, with user-memcpy

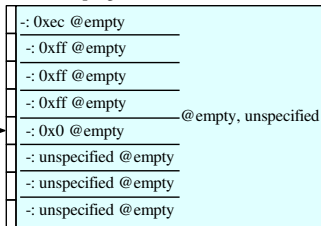
n: size_t [@7, 0xfffffc0]



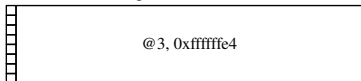
dest: unsigned char* [@5, 0xfffffd0]



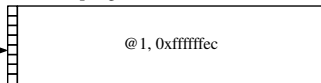
q: signed int* [@4, 0xfffffd8]



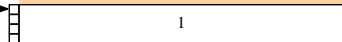
src: unsigned char* [@6, 0xfffffc8]



p: signed int* [@3, 0xfffffe0]

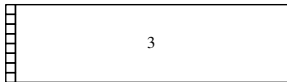


x: signed int [@1 exp, 0xfffffec]

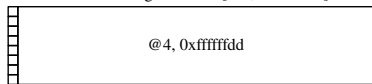


Copying pointer values byte-wise, with user-memcpy

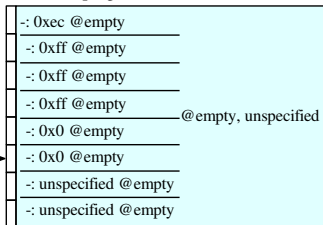
n: size_t [@7, 0xfffffc0]



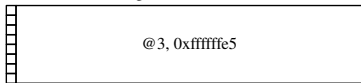
dest: unsigned char* [@5, 0xfffffd0]



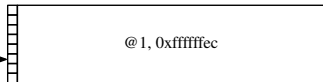
q: signed int* [@4, 0xfffffd8]



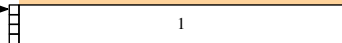
src: unsigned char* [@6, 0xfffffc8]



p: signed int* [@3, 0xfffffe0]

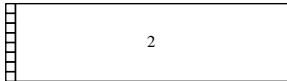


x: signed int [@1 exp, 0xfffffec]

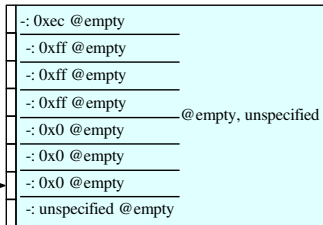


Copying pointer values bitwise, with user-memcpy

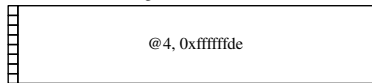
n: size_t [@7, 0xffffffc0]



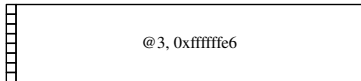
q: signed int* [@4, 0xfffffd8]



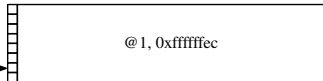
dest: unsigned char* [@5, 0xfffffd0]



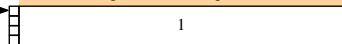
src: unsigned char* [@6, 0xfffffc8]



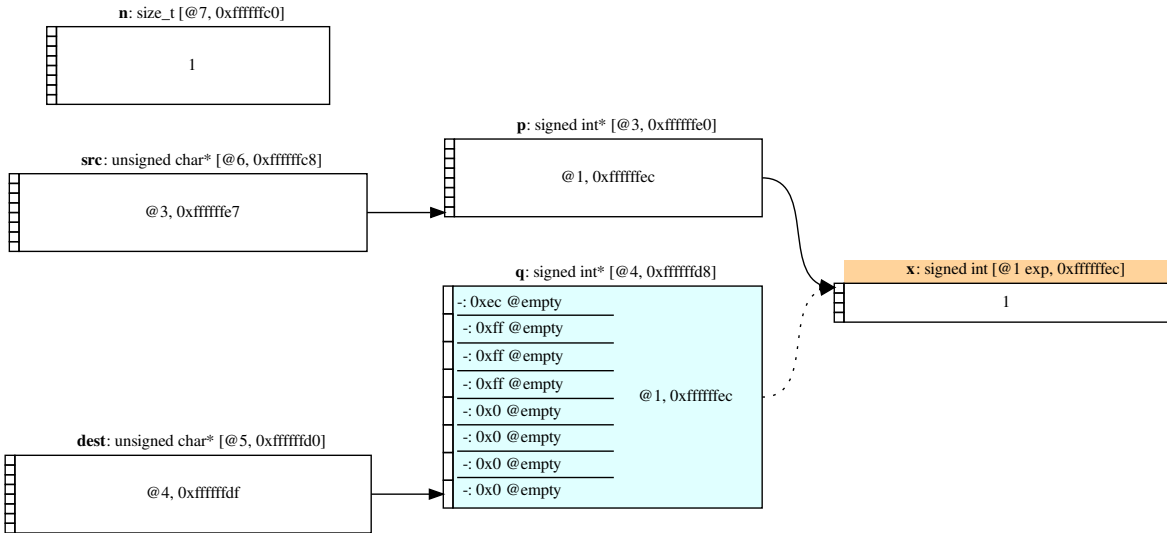
p: signed int* [@3, 0xfffffe0]



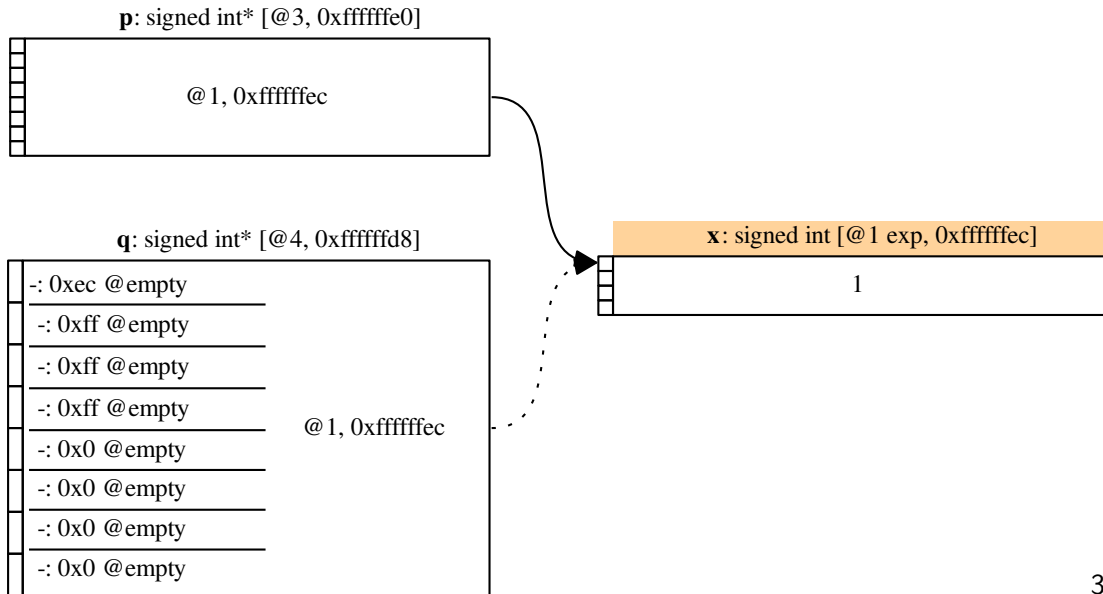
x: signed int [@1 exp, 0xfffffec]



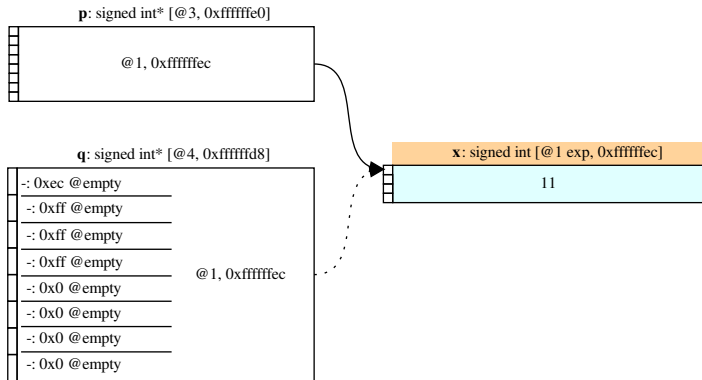
Copying pointer values bitwise, with user-memcpy



Copying pointer values bitwise, with user-memcpy



Copying pointer values bitwise, with user-memcpy



The first read of a `p` pointer byte marked `x` as exposed, then the final `*q=11` access follows the integer-to-pointer cast semantics when reading a pointer value from the memory bytes, recovering the provenance `@1` that the concrete address is within.

Pointer provenance and union type punning

Pointer values can also be constructed by type punning, e.g. writing an `int*` union member, reading it as a `uintptr_t` union member, and then casting back to a pointer type.

(The example assumes the object representations of the pointer and the result of the cast to integer are identical. This is not guaranteed by the standard, but holds for many implementations.)

```
// provenance_union_punning_3_global.c
1 #include <stdio.h>
2 #include <string.h>
3 #include <inttypes.h>
4 int x=1;
5 typedef union { uintptr_t ui; int *up; } un;
6 int main() {
7     un u;
8     int *p = &x;
9     u.up = p;
10    uintptr_t i = u.ui;
11    int *q = (int*)i;
12    *q = 11; // does this have UB?
13    printf("x=%d *p=%d *q=%d\n", x, *p, *q);
14    return 0;
15 }
```

The same semantics as for representation-byte reads also permits this: `x` is deemed exposed by the read of the provenanced representation bytes by the non-pointer-type read. The integer-to-pointer cast then recreates the provenance of `x`.

Pointer provenance via IO

Three versions:

- ▶ using `fprintf/fscanf` and the `%p` format (which the standard says should work),
- ▶ using `fwrite/fread` on the pointer representation bytes, and
- ▶ converting the pointer to and from `uintptr_t` and using `fprintf/fscanf`.

The first gives a syntactic indication of a potentially escaping pointer value; the others do not.

Exotic, but used in practice.

In our proposal, these just work: we mark the storage instance as exposed on the `%p` printf, pointer representation-byte read, or cast, and use the same semantics as integer-to-pointer casts at input-, read-, or cast-time to recover the original provenance.

Implications for optimisation

Can a function argument alias its local variables? (1/3)

This should be forbidden:

```
// pointer_from_integer_1pg.c
1 #include <stdio.h>
2 #include <stdint.h>
3 #include "charon_address_guesses.h"
4 void f(int *p) {
5     int j=5;
6     if (p==&j)
7         *p=7;
8     printf("j=%d &j=%p\n", j, (void*)&j);
9 }
10 int main() {
11     uintptr_t i = ADDRESS_PFI_1PG;
12     int *p = (int*)i;
13     f(p);
14 }
```

main() guesses the address of f()'s local variable j, passing it in as a pointer, and f() checks it before using it for an access. Here GCC -O0 optimises away the `if` and the write `*p=7`, even when `ADDRESS_PFI_1PG` is the same as `&j`. That compiler behaviour should be permitted, so this program should be deemed UB. In other words, code should not normally be allowed to rely on implementation facts about the allocation addresses of C variables.

Can a function argument alias its local variables? (1/3)

This should be forbidden:

```
// pointer_from_integer_1pg.c
1 #include <stdio.h>
2 #include <stdint.h>
3 #include "charon_address_guesses.h"
4 void f(int *p) {
5     int j=5;
6     if (p==&j)
7         *p=7;
8     printf("j=%d &j=%p\n", j, (void*)&j);
9 }
10 int main() {
11     uintptr_t i = ADDRESS_PFI_1PG;
12     int *p = (int*)i;
13     f(p);
14 }
```

Our PNVI-* proposals correctly deems this to be UB: at the point of the `(int*)i` cast the `j` storage instance does not yet even exist, so that cast gives a pointer with empty provenance; any execution that goes into the `if` would thus flag UB, so the program as a whole is UB.

Can a function argument alias its local variables? (2/3)

Varying to do the (`int*`) cast after the `j` allocation:

```
// pointer_from_integer_lig.c
1 #include <stdio.h>
2 #include <stdint.h>
3 #include "charon_address_guesses.h"
4 void f(uintptr_t i) {
5     int j=5;
6     int *p = (int*)i;
7     if (p==&j)
8         *p=7;
9     printf("j=%d &j=%p\n", j, (void*)&j);
10 }
11 int main() {
12     uintptr_t j = ADDRESS_PFI_1IG;
13     f(j);
14 }
```

This is still forbidden in PNVI-ae-*, as `j` is not exposed. It would be allowed in PNVI-plain, but perhaps that would also be acceptable – it would just require compilers to be conservative about the results of integer-to-pointer casts where they cannot see the source of the integer, which we imagine is a rare case.

Can a function argument alias its local variables? (3/3)

Varying again to remove the conditional guard and make `j` exposed:

```
// pointer_from_integer_lie.c
1 #include <stdio.h>
2 #include <stdint.h>
3 #include "charon_address_guesses.h"
4 void f(uintptr_t i) {
5     int j=5;
6     uintptr_t k = (uintptr_t)&j;
7     int *p = (int*)i;
8     *p=7;
9     printf("j=%d\n",j);
10 }
11 int main() {
12     uintptr_t j = ADDRESS_PFI_1I;
13     f(j);
14 }
```

Executions in which `&j == ADDRESS_PFI_1I` would be ok, but, because the standard does not and should not constrain allocation addresses (beyond alignment and non-overlapping properties), there are always (unless the address space is almost exhausted) other executions in which `ADDRESS_PFI_1I` does not match any allocation. So this is still (correctly) deemed UB.

Can a function argument alias its local variables? (3/3)

Varying again to remove the conditional guard and make `j` exposed:

```
// pointer_from_integer_lie.c
1 #include <stdio.h>
2 #include <stdint.h>
3 #include "charon_address_guesses.h"
4 void f(uintptr_t i) {
5     int j=5;
6     uintptr_t k = (uintptr_t)&j;
7     int *p = (int*)i;
8     *p=7;
9     printf("j=%d\n",j);
10 }
11 int main() {
12     uintptr_t j = ADDRESS_PFI_1I;
13     f(j);
14 }
```

In other words: the fact that programmers cannot assume anything about allocation addresses licenses the desired compiler optimisation. That's expressed in the abstract machine simply by making allocation addresses nondeterministic.

Can a function access local variables of its parent? (1/2)

This too should be forbidden in general.

```
// pointer_from_integer_2.c
1 #include <stdio.h>
2 #include <stdint.h>
3 #include "charon_address_guesses.h"
4 void f() {
5     uintptr_t i=ADDRESS_PFI_2;
6     int *p = (int*)i;
7     *p=7;
8 }
9 int main() {
10    int j=5;
11    f();
12    printf("j=%d\n", j);
13 }
```

Here `f()` guesses the address of `main()`'s local variable `j`.

This is similarly UB by allocation-address nondeterminism: the abstract machine allows executions in which the guess is correct, but also executions in which it is incorrect, where the `*p=7` flags UB. So the program is UB.

Can a function access local variables of its parent? (1/2)

This too should be forbidden in general.

```
// pointer_from_integer_2.c
1 #include <stdio.h>
2 #include <stdint.h>
3 #include "charon_address_guesses.h"
4 void f() {
5     uintptr_t i=ADDRESS_PFI_2;
6     int *p = (int*)i;
7     *p=7;
8 }
9 int main() {
10    int j=5;
11    f();
12    printf("j=%d\n", j);
13 }
```

Here `f()` guesses the address of `main()`'s local variable `j`.

(In PNVI-ae-*, `j` is not exposed, so all executions flag UB, but the previous argument applies even if `j` is exposed.)

Can a function access local variables of its parent? (2/2)

Varying to guard the call to `f()` with an address check:

```
// pointer_from_integer_2g.c
1 #include <stdio.h>
2 #include <stdint.h>
3 #include "charon_address_guesses.h"
4 void f() {
5     uintptr_t i=ADDRESS_PFI_2G;
6     int *p = (int*)i;
7     *p=7;
8 }
9 int main() {
10    int j=5;
11    if ((uintptr_t)&j == ADDRESS_PFI_2G)
12        f();
13    printf("j=%d &j=%p\n", j, (void*)&j);
14 }
```

This is allowed in PNV1-*, but the guard necessarily involves `&j`, so compilers should be able to deem this escaped. In other words, while we don't think this example needs to be allowed, it should be ok to make it allowed.

Optimisations based on equality tests

In any provenance-aware semantics, $p==q$ can hold in some cases where p and q are not interchangeable (e.g. $*p$ is defined but $*q$ UB).

(Otherwise, we'd have to require implementations track provenance at runtime for $==$ testing; not usually practical.)

As Lee et al. observe [OOPSLA 2018], that restricts optimisations, e.g. GVN, based on pointer equality tests.

Solution: just don't do those.

(There's no alternative, short of compilers giving up on provenance-based alias analysis altogether, which would be worse.)

PNVI-plain vs PNVI-ae-*

Is the PNVI-ae-* “exposed” machinery necessary?

Debatable. There’s not much difference between PNVI-plain and PNVI-ae for these examples ([pointer_from_integer_1ig.c](#) is allowed in PNVI-plain but forbidden in PNVI-ae-*).

PNVI-plain is simpler, but relies on allocation-address nondeterminism (which some people aren’t happy with) for more of the examples than PNVI-ae-*.

PNVI-ae-* is more complex, but makes some of these examples UB just by examining a single execution path. It’s also subject to...

The problem with lost address-takens and escapes

```
// provenance_lost_escape_1.c
1 #include <stdio.h>
2 #include <string.h>
3 #include <stdint.h>
4 #include "charon_address_guesses.h"
5 int x=1; // assume allocation ID @1, at ADDR_PLE_1
6 int main() {
7     int *p = &x;
8     uintptr_t i1 = (intptr_t)p;           // (@1,ADDR_PLE_1)
9     uintptr_t i2 = i1 & 0x00000000FFFFFFFF; //
10    uintptr_t i3 = i2 & 0xFFFFFFFF00000000; // (@1,0x0)
11    uintptr_t i4 = i3 + ADDR_PLE_1;       // (@1,ADDR_PLE_1)
12    int *q = (int *)i4;
13    printf("Addresses: p=%p\n", (void*)p);
14    if (memcmp(&i1, &i4, sizeof(i1)) == 0) {
15        *q = 11; // does this have defined behaviour?
16        printf("x=%d *p=%d *q=%d\n", x, *p, *q);
17    }
18 }
```

In PNVI-plain, this is allowed, simply because `x` exists at the integer-to-pointer cast. Implementations that are conservative w.r.t. all pointers formed from integers would automatically be sound w.r.t. that.

The problem with lost address-takens and escapes

```
// provenance_lost_escape_1.c
1 #include <stdio.h>
2 #include <string.h>
3 #include <stdint.h>
4 #include "charon_address_guesses.h"
5 int x=1; // assume allocation ID @1, at ADDR_PLE_1
6 int main() {
7     int *p = &x;
8     uintptr_t i1 = (intptr_t)p;           // (@1,ADDR_PLE_1)
9     uintptr_t i2 = i1 & 0x00000000FFFFFFFF; //
10    uintptr_t i3 = i2 & 0xFFFFFFFF00000000; // (@1,0x0)
11    uintptr_t i4 = i3 + ADDR_PLE_1;       // (@1,ADDR_PLE_1)
12    int *q = (int *)i4;
13    printf("Addresses: p=%p\n", (void*)p);
14    if (memcmp(&i1, &i4, sizeof(i1)) == 0) {
15        *q = 11; // does this have defined behaviour?
16        printf("x=%d *p=%d *q=%d\n", x, *p, *q);
17    }
18 }
```

In PNVI-ae-*, in the source program x is exposed before the integer-to-pointer cast, so this is allowed here too.

But a compiler might optimise (in its intermediate language)...

The problem with lost address-takens and escapes

```
// provenance_lost_escape_1_optimised.c
1 #include <stdio.h>
2 #include <string.h>
3 #include <stdint.h>
4 #include "charon_address_guesses.h"
5 int x=1; // assume allocation ID @1, at ADDR_PLE_1
6 int main() {
7     int *p = &x;
8
9
10
11     uintptr_t i4 = ADDR_PLE_1;
12     int *q = (int *)i4;
13     printf("Addresses: p=%p\n", (void*)p);
14     uintptr_t i1 = (uintptr_t)p;
15     if (memcmp(&i1, &i4, sizeof(i1)) == 0) {
16         *q = 11; // does this have defined behaviour?
17         printf("x=%d *p=%d *q=%d\n", x, *p, *q);
18     }
19 }
```

and now `x` is no longer exposed before the cast. If this happens before alias analysis, the results would be wrong.

The problem with lost address-takens and escapes

Solutions: either

- ▶ simply be conservative (in alias analysis) w.r.t. all pointers formed from integers, or
- ▶ record, in optimisations that occur before alias analysis, any lost exposures, and pass those in as an additional argument to alias analysis.

PNVI-ae vs PNVI-ae-udi

Should we allow one-past integer-to-pointer casts?

We have to decide whether casting a one-past pointer to integer and back gives a usable result.

```
// provenance_roundtrip_via_intptr_t_onepast.c
1 #include <stdio.h>
2 #include <inttypes.h>
3 int x=1;
4 int main() {
5     int *p = &x;
6     p=p+1;
7     intptr_t i = (intptr_t)p;
8     int *q = (int *)i;
9     q=q-1;
10    *q = 11; // is this free of undefined behaviour?
11    printf("*p=%d *q=%d\n",*p,*q);
12 }
```

Pro: it's nice for one-past pointers to behave like in-bounds pointers

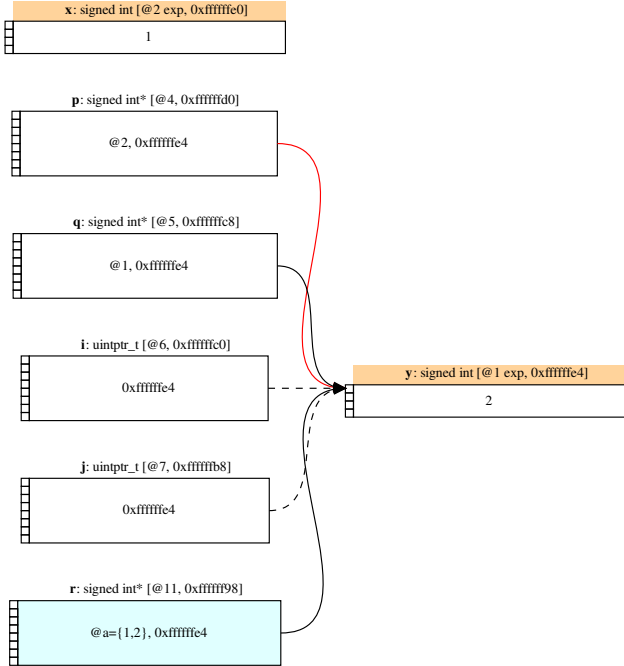
Con: if that's allowed, we have to deal with ambiguous integers, which can be regarded either one-past one object or the start of another.

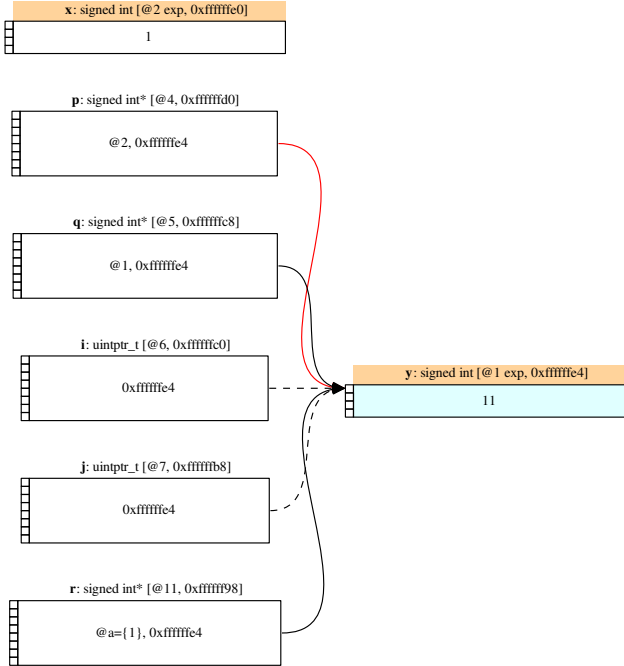
Should we allow one-past integer-to-pointer casts?

We have to decide whether casting a one-past pointer to integer and back gives a usable result.

PNVI-plain and PNVI-ae forbid this: an integer has to be properly within an object for it to be castable to a usable pointer.

PNVI-ae-udi (user disambiguation) permits it: it leaves the provenance of pointer values resulting from such casts unknown until the first operation (e.g. an access, pointer arithmetic, or pointer relational comparison) that disambiguates them. This makes examples that use the result of the cast in one consistent way well defined.





Experimental checks

Testing the example behaviour in Cerberus

We confirmed the examples behave as desired in each model by running them in Cerberus.

test family	test	intended behaviour			observed behaviour	
		PNVI-plain	PNVI-aa	PNVI-aa-udi	PNVI-plain Cerberus (decreasing allocator)	PNVI-aa PNVI-aa-udi
1	provenance_basic_global_xy.c				not triggered	
	provenance_basic_global_xy.c		UB		UB (line 9)	
	provenance_basic_auto_xy.c				not triggered	
2	cheri_03_i.c			UB	UB (except with <code>permissive_pointer_arith</code> switch)	
	pointer_offset_from_ptr_subtraction_global_xy.c				UB (pointer subtraction)	
	pointer_offset_from_ptr_subtraction_auto_xy.c		UB (pointer subtraction)		UB (out-of-bound store with <code>permissive_pointer_arith</code> switch)	
3	provenance_equality_global_xy.c				not triggered	
	provenance_equality_global_xy.c				defined (ND except with <code>strict_pointer_equality</code> switch)	
	provenance_equality_auto_xy.c				not triggered	
	provenance_equality_auto_xy.c		defined, nondet		defined (ND except with <code>strict_pointer_equality</code> switch)	
	provenance_equality_global_in_xy.c				not triggered	
4	provenance_equality_global_in_xy.c				defined (ND except with <code>strict_pointer_equality</code> switch)	
	provenance_roundtrip_via_intptr_1.c		defined		defined	
5	provenance_basic_using_uintptr_1_global_xy.c				not triggered	
	provenance_basic_using_uintptr_1_global_xy.c		defined		defined	
	provenance_basic_using_uintptr_1_auto_xy.c				not triggered	
	provenance_basic_using_uintptr_1_auto_xy.c				defined	
6	pointer_offset_from_int_subtraction_global_xy.c				defined	
	pointer_offset_from_int_subtraction_global_xy.c		defined		defined	
	pointer_offset_from_int_subtraction_auto_xy.c				defined	
	pointer_offset_from_int_subtraction_auto_xy.c				defined	
7	pointer_offset_xor_global.c				defined	
	pointer_offset_xor_auto.c		defined		defined	
8	pointer_offset_xor_global.c				defined	
	pointer_offset_xor_auto.c		defined		defined	
9	provenance_tag_bits_via_uintptr_1_1.c		defined		defined	
	pointer_arith_algebraic_properties_2_global.c		defined		defined	
	pointer_arith_algebraic_properties_3_global.c		defined		defined	
	pointer_copy_memory.c		defined		defined	
	pointer_copy_user_defined_direct_bytewise.c		defined		defined	
	provenance_tag_bits_via_repr_byte_1.c		defined		defined	
	pointer_copy_user_defined_bytewise.c		defined		defined	
	pointer_copy_user_defined_bytewise.c		defined		defined	
	provenance_equality_uintptr_1_global_xy.c				not triggered	
	provenance_equality_uintptr_1_global_xy.c		defined		defined (true)	
	provenance_equality_uintptr_1_auto_xy.c				not triggered	
10	provenance_equality_uintptr_1_auto_xy.c				defined (true)	
	provenance_union_punning_2_global_xy.c	defined	UB (line 16, deref)	UB (line 16, store)	not triggered	
	provenance_union_punning_2_global_xy.c	defined	UB (line 16, deref)	UB (line 16, store)	defined	UB (line 16, deref) UB (line 16, store)
	provenance_union_punning_2_auto_xy.c	defined	UB (line 16, deref)	UB (line 16, store)	not triggered	
	provenance_union_punning_2_auto_xy.c	defined	UB (line 16, deref)	UB (line 16, store)	defined	UB (line 16, deref) UB (line 16, store)
11	provenance_union_punning_3_global.c				defined	
	provenance_via_io_percentp_global.c				defined	
	provenance_via_io_bytewise_global.c				defined	
	provenance_via_io_uintptr_1_global.c				defined	
	pointer_from_integer_1p.c		UB (line 7)		UB (line 7)	UB (line 7)
	pointer_from_integer_1p.c	defined (j = 7)	UB (line 8)		defined (j = 7)	UB (line 8)
	pointer_from_integer_1c.c		UB (line 6)		UB (line 6)	UB (line 6)
	pointer_from_integer_1c.c	defined (j = 7)	UB (line 7)		defined (j = 7)	UB (line 7)
	pointer_from_integer_2c.c		defined (j = 7)		defined (j = 7)	UB (line 7)
	pointer_from_integer_2c.c	defined (j = 7)	UB (line 7)		defined (j = 7)	UB (line 7)
	pointer_from_integer_2g.c		defined (j = 7)		defined (j = 7)	UB (line 7)
	provenance_test_ctype_1.c		defined		defined	defined
12	provenance_roundtrip_via_intptr_1_onepass.c	UB (line 10)	defined	defined	UB (line 10)	defined
	pointer_from_int_disambiguation_1c.c		defined (y = 11)		defined (y = 11)	
13	pointer_from_int_disambiguation_1_xy.c				not triggered	
	pointer_from_int_disambiguation_2c.c				UB (line 14)	defined (x = 11)
	pointer_from_int_disambiguation_2_xy.c	UB (line 14)	defined		not triggered	
	pointer_from_int_disambiguation_3c.c				UB (line 15)	
14	pointer_from_int_disambiguation_3_xy.c				not triggered	

(bold = tests mentioned in the document)

green = Cerberus behaviour matches intent

blue = Cerberus behaviour matches intent (with `permissive_pointer_arith` switch)

grey = Cerberus' allocator doesn't trigger the interesting behaviour

Testing the example behaviour in mainstream C implementations

It doesn't seem possible to make a coherent and useful semantics that admits all the existing observed compiler behaviour – but they do agree in many cases, and it may be that only mild adaptations would be needed.

Pointer equality

Consider pointers p and q with different provenance. In an execution where they have the same address (same pointer object representation), is $p==q$:

1. required to be true, or
2. allowed to be either true or false, or
3. undefined behaviour?

C17 6.5.9p6 says (1) *“Two pointers compare equal **if and only if** both are [...] or one is a pointer to one past the end of one array object and the other is a pointer to the start of a different array object that happens to immediately follow the first array object in the address space”*,

GCC follows (2).

We suspect (3) would break existing code.

Pick one...