

# C provenance semantics: examples

Peter Sewell, Kayvan Memarian, Victor B F Gomes, Jens Gustedt, Martin Uecker

## ▶ To cite this version:

Peter Sewell, Kayvan Memarian, Victor B F Gomes, Jens Gustedt, Martin Uecker. C provenance semantics: examples. [Technical Report] N2363, ISO JCT1/SC22/WG14. 2019. hal-02089907

# HAL Id: hal-02089907 https://inria.hal.science/hal-02089907

Submitted on 4 Apr 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## C provenance semantics: examples

2 (for PNVI-plain, PNVI address-exposed, PNVI address-exposed user-disambiguation, and PVI models) 3 4 PETER SEWELL, University of Cambridge 5 KAYVAN MEMARIAN, University of Cambridge 6 VICTOR B. F. GOMES, University of Cambridge 7 JENS GUSTEDT, Université de Strasbourg, CNRS, Inria, ICube, Strasbourg, France 8 MARTIN UECKER, University Medical Center, Goettingen 9 10 This note discusses the design of provenance semantics for C, looking at a series of examples. We consider three variants 11 of the provenance-not-via-integer (PNVI) model: PNVI plain, PNVI address-exposed (PNVI-ae) and PNVI address-exposed 12 user-disambiguation (PNVI-ae-udi), and also the provenance-via-integers (PVI) model. The examples include those of Exploring 13 C Semantics and Pointer Provenance [POPL 2019] (also available as ISO WG14 N2311 http://www.open-std.org/jtc1/sc22/wg14/ www/docs/n2311.pdf), with several additions. This is based on recent discussion in the C memory object model study group. It 14 should be read together with the two companion notes, one giving detailed diffs to the C standard text (N2362), and another 15 giving detailed semantics for these variants (N2364). 16 Contents 17 Abstract 1 18 19 Contents 1 20 1 Introduction 1 2 Basic pointer provenance 2 21 22 3 Refining the basic provenance model to support pointer construction via casts, representation 5 23 accesses, etc. 24 4 Refining the basic provenance model: phenomena and examples 6 25 5 Implications of provenance semantics for pptimisations 13 Testing the example behaviour in Cerberus 26 6 19 27 7 Testing the example behaviour in mainstream C implementations 20 28 References 21 29

### 1 INTRODUCTION

30

31

32

1

The new material for PNVI-address-exposed and PNVI address-exposed user-disambiguation models starts in §3, but first we introduce introduce the problem in general and describe the basic pointer provenance semantics.

The semantics of pointers and memory objects in C has been a vexed question for many years. A priori, 33 one might imagine two language-design extremes: a concrete model that exposes the memory semantics of the 34 underlying hardware, with memory being simply a finite partial map from machine-word addresses to bytes 35 and pointers that are simply machine words, and an abstract model in which the language types enforce hard 36 distinctions, e.g. between numeric types that support arithmetic and pointer types that support dereferencing. C 37 is neither of these. Its values are not abstract: the language intentionally permits manipulation of their underlying 38 representations, via casts between pointer and integer types, char\* pointers to access representation bytes, and 39 so on, to support low-level systems programming. But C values also cannot be considered to be simple concrete 40 values: at runtime a C pointer will typically just be a machine word, but compiler analysis reasons about abstract 41 notions of the provenance of pointers, and compiler optimisations rely on assumptions about these for soundness. 42 Particularly relevant here, some compiler optimisations rely on alias analysis to deduce that two pointer values do 43 not refer to the same object, which in turn relies on assumptions that the program only constructs pointer values 44 in "reasonable" ways (with other programs regarded as having undefined behaviour, UB). The committee response 45 to Defect Report DR260 [Feather 2004] states that implementations can track the origins (or "provenance") of 46 pointer values, "the implementation is entitled to take account of the provenance of a pointer value when determining 47 what actions are and are not defined", but exactly what this "provenance" means is left undefined, and it has never 48 been incorporated into the standard text. Even what a memory object is is not completely clear in the standard, 49 especially for aggregate types and for objects within heap regions. 50

Second, in some respects there are significant discrepancies between the ISO standard and the de facto standards, of C as it is implemented and used in practice. Major C codebases typically rely on particular compiler flags, e.g. -fno-strict-aliasing or -fwrapv, that substantially affect the semantics but which standard does not attempt to describe, and some idioms are UB in ISO C but relied on in practice, e.g. comparing against a pointer value after the lifetime-end of the object it pointed to. There is also not a unique de facto standard: in reality, one has to consider the expectations of expert C programmers and compiler writers, the behaviours of specific compilers, and the assumptions about the language implementations that the global C codebase relies upon to

2019.

59 60

work correctly (in so far as it does). Our recent surveys [Memarian et al. 2016; Memarian and Sewell 2016b] of the first revealed many discrepancies, with widely conflicting responses to specific questions.

Third, the ISO standard is a prose document, as is typical for industry standards. The lack of mathematical 63 precision, while also typical for industry standards, has surely contributed to the accumulated confusion about C, 64 but, perhaps more importantly, the prose standard is not executable as a test oracle. One would like, given small 65 test programs, to be able to automatically compute the sets of their allowed behaviours (including whether they 66 have UB). Instead, one has to do painstaking argument with respect to the text and concepts of the standard, a 67 time-consuming and error-prone task that requires great expertise, and which will sometimes run up against the 68 areas where the standard is unclear or differs with practice. One also cannot use conventional implementations to 69 find the sets of all allowed behaviours, as (a) the standard is a loose specification, while particular compilations 70 will resolve many nondeterministic choices, and (b) conventional implementations cannot detect all sources of 71 undefined behaviour (that is the main point of UB in the standard, to let implementations assume that source 72 programs do not exhibit UB, together with supporting implementation variation beyond the UB boundary). 73 Sanitisers and other tools can detect some UB cases, but not all, and each tool builds in its own more-or-less ad 74 hoc C semantics. 75

This is not just an academic problem: disagreements over exactly what is or should be permitted in C have caused considerable tensions, e.g. between OS kernel and compiler developers, as increasingly aggressive optimisations can break code that worked on earlier compiler implementations.

This note continues an exploration of the design space and two candidate semantics for pointers and memory 79 objects in C, taking both ISO and de facto C into account. We earlier [Chisnall et al. 2016; Memarian et al. 2016] 80 identified many design questions. We focus here on the questions concerning pointer provenance, which we 81 revise and extend. We develop two main coherent proposals that reconcile many design concerns; both are 82 broadly consistent with the provenance intuitions of practitioners and ISO DR260, while still reasonably simple. 83 We highlight their pros and cons and various outstanding open questions. These proposals cover many of the 84 interactions between abstract and concrete views in C: casts between pointers and integers, access to the byte 85 representations of values, etc. 86

### 2 BASIC POINTER PROVENANCE

<sup>89</sup> C pointer values are typically represented at runtime as simple concrete numeric values, but mainstream compilers
 <sup>90</sup> routinely exploit information about the *provenance* of pointers to reason that they cannot alias, and hence to
 <sup>91</sup> justify optimisations. In this section we develop a provenance semantics for simple cases of the construction and
 <sup>92</sup> use of pointers,

93 For example, consider the classic 94 test [Chisnall et al. 2016; Feather 2004; 95 Krebbers and Wiedijk 2012; Krebbers 2015; 96 Memarian et al. 2016] on the right (note 97 that this and many of the examples below 98 are edge-cases, exploring the boundaries 99 of what different semantic choices allow, 100 and sometimes what behaviour existing 101 compilers exhibit; they are not all intended 102 as desirable code idioms).

Depending on the implementation, x and
 y might in some executions happen to be al located in adjacent memory, in which case

```
// provenance_basic_global_yx.c (and an xy variant)
#include <stdio.h>
#include <string.h>
int y=2, x=1;
int main() {
    int *p = &x + 1;
    int *q = &y;
    printf("Addresses: p=%p q=%p\n",(void*)p,(void*)q);
    if (memcmp(&p, &q, sizeof(p)) == 0) {
        *p = 11; // does this have undefined behaviour?
        printf("x=%d y=%d *p=%d *q=%d\n",x,y,*p,*q);
    }
}
```

106 &x+1 and &y will have bitwise-identical representation values, the memcmp will succeed, and p (derived from a pointer 107 to x) will have the same representation value as a pointer to a different object, y, at the point of the update \*p=11. 108 This can occur in practice, e.g. with GCC 8.1 -O2 on some platforms. Its output of x=1 y=2 \*p=11 \*q=2 suggests 109 that the compiler is reasoning that \*p does not alias with y or \*q, and hence that the initial value of y=2 can be 110 propagated to the final printf. ICC, e.g. ICC 19 -O2, also optimises here (for a variant with x and y swapped), 111 producing x=1 y=2 \*p=11 \*q=11. In contrast, Clang 6.0 -O2 just outputs the x=1 y=11 \*p=11 \*q=11 that one 112 might expect from a concrete semantics. Note that this example does not involve type-based alias analysis, and 113 the outcome is not affected by GCC or ICC's -fno-strict-aliasing flag. Note also that the mere formation of the 114 &x+1 one-past pointer is explicitly permitted by the ISO standard, and, because the \*p=11 access is guarded by the 115 memcmp conditional check on the representation bytes of the pointer, it will not be attempted (and hence flag UB) 116 in executions in which the two storage instances are not adjacent.

These GCC and ICC outcomes would not be correct with respect to a concrete semantics, and so to make the existing compiler behaviour sound it is necessary for this program to be deemed to have undefined behaviour.

61

62

87

88

125

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

146

147

148

The current ISO standard text does not explicitly speak to this, but the 2004 ISO WG14 C standards committee 121 response to Defect Report 260 (DR260 CR) [Feather 2004] hints at a notion of provenance associated to values that 122 keeps track of their "origins": 123

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

126 However, DR260 CR has never been incorporated in the standard text, and it gives no more detail. This leaves 127 many specific questions unclear: it is ambiguous whether some programming idioms are allowed or not, and 128 exactly what compiler alias analysis and optimisation are allowed to do.

**Basic provenance semantics for pointer values** For simple cases of the construction and use of pointers, capturing the basic intuition suggested by DR260 CR in a precise semantics is straightforward: we associate a provenance with every pointer value, identifying the original storage instance the pointer is derived from. In more detail:

- We take abstract-machine pointer values to be pairs  $(\pi, a)$ , adding a *provenance*  $\pi$ , either @*i* where *i* is a storage instance ID, or the empty provenance @empty, to their concrete address a.
- On every storage instance (of objects with static, thread, automatic, and allocated storage duration), the abstract machine nondeterministically chooses a fresh storage instance ID i (unique across the entire execution), and the resulting pointer value carries that single storage instance ID 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. In particular:
- access via a pointer value which has provenance a single storage instance ID @i must be within the memory footprint of the corresponding original storage instance, which must still be live.
- all other accesses, including those via a pointer value with empty provenance, are undefined behaviour. Regarding such accesses as undefined behaviour is necessary to make optimisation based on provenance 145 alias analysis sound: if the standard did define behaviour for programs that make provenance-violating accesses, e.g. by adopting a concrete semantics, optimisation based on provenance-aware alias analysis would not be sound.

On the right is a provenance-semantics memory-149 state snapshot (from the Cerberus GUI) for 150 provenance\_basic\_global\_xy.c, just before the 151 152 invalid access via p, showing how the provenance mismatch makes it UB: at the attempted access via p, 153 its pointer-value address 0x4c is not within the storage 154 instance with the ID @5 of the provenance of p. 155

All this is for the *C* abstract machine as defined in 156 the standard: compilers might rely on provenance in 157 their alias analysis and optimisation, but one would 158 not expect normal implementations to record or ma-159 nipulate provenance at runtime (though dynamic or 160



static analysis tools might, as might non-standard implementations such as CHERI C). Provenances therefore do 161 not have program-accessible runtime representations in the abstract machine. 162

Even for the basic provenance semantics, there are some open design questions, which we now discuss. 163

164 Can one construct out-of-bounds (by more than one) pointer values by pointer arithmetic? Consider 165 the example below, where q is transiently (more than one-past) out of bounds but brought back into bounds before 166 being used for access. In ISO C, constructing such a pointer value is clearly stated to be undefined behaviour [WG14 167 2018, 6.5.6p8]. This can be captured using the provenance of the pointer value to determine the relevant bounds. 168 There are cases where such pointer arithmetic would go

169 wrong on some platforms (some now exotic), e.g. where 170 pointer arithmetic subtraction overflows, or if the transient 171 value is not aligned and only aligned values are representable 172 at the particular pointer type, or for hardware that does 173 bounds checking, or where pointer arithmetic might wrap

174 at values less than the obvious word size (e.g. "near" or "huge" 8086 pointers). However, transiently out-of-bounds 175 pointer construction seems to be common in practice. It may be desirable to make it implementation-defined 176 whether such pointer construction is allowed. That would continue to permit implementations in which it would 177 go wrong to forbid it, but give a clear way for other implementations to document that they do not exploit this 178 UB in compiler optimisations that may be surprising to programmers. 179

**Inter-object pointer arithmetic** The first example in this section relied on guessing (and then checking) the 181 offset between two storage instances. What if one instead calculates the offset, with pointer subtraction; should that 182 let one move between objects, as below? In ISO C18, the q-p is UB (as it is a pointer subtraction between pointers to 183

#include <stddef.h>

**int** x=1, y=2;

int \*p = &x;

**int** \*q = &y;

ptrdiff\_t offset = q - p;

if (memcmp(&r, &q, sizeof(r)) == 0) {

printf("y=%d \*q=%d \*r=%d\n",y,\*q,\*r);

\*r = 11; // is this free of UB?

int \*r = p + offset;

int main() {

}

}

```
different objects, which in some abstract-
184
      machine executions are not one-past-related). In
185
                                                       #include <stdio.h>
      a variant semantics that allows construction of
                                                       #include <string.h>
186
```

more-than-one-past pointers (which allows the 187

evaluation of p + offset), one would have to to 188

choose whether the \*r=11 access is UB or not. 189

The basic provenance semantics will forbid it, 190

because r will retain the provenance of the x stor-191 age instance, but its address is not in bounds for

192 that. This is probably the most desirable seman-

193 tics: we have found very few example idioms that 194

intentionally use inter-object pointer arithmetic, 195

- and the freedom that forbidding it gives to alias 196
- analysis and optimisation seems significant. 197

```
198
199
```

200 201

```
Pointer equality comparison and provenance
```

202 203 with the same address but different provenance 204 as nonequal. Unsurprisingly, this happens in 205 some circumstances but not others, e.g. if the 206 test is pulled into a simple separate function, but 207 not if in a separate compilation unit. To be conser-

208 vative w.r.t. current compiler behaviour, pointer

209 equality in the semantics should give false if the

210 addresses are not equal, but nondeterministically 211

(at each runtime occurrence) either take prove-212 nance into account or not if the addresses are

```
213
      equal - this specification looseness accommo-
```

```
214
      dating implementation variation. Alternatively,
```

```
215
      one could require numeric comparisons, which
```

216 would be a simpler semantics for programmers

A priori, pointer equality comparison (with == or !=) might be expected to just compare their numeric addresses, but we observe GCC 8.1 -O2 sometimes regarding two pointers

// pointer\_offset\_from\_ptr\_subtraction\_global\_xy.c

```
// provenance_equality_global_xy.c
#include <stdio.h>
#include <string.h>
int x=1, y=2;
int main() {
  int *p = &x + 1;
  int *q = &y;
  printf("Addresses: p=%p q=%p\n",(void*)p,(void*)q);
  Bool b = (p==a):
  // can this be false even with identical addresses?
  printf("(p==q) = %s\n", b?"true":"false");
  return 0;
```

217 but force that GCC behaviour to be regarded as a bug. Cerberus supports both options. One might also imagine 218 making it UB to compare pointers that are not strictly within their original storage instance [Krebbers 2015], but 219 that would break loops that test against a one-past pointer, or requiring equality to *always* take provenance into 220 account, but that would require implementations to track provenance at runtime.

}

The current ISO C18 standard text is too strong here unless numeric comparison is required: 6.5.9p6 says "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", which requires such pointers to compare equal – reasonable pre-DR260 CR, but debatable after it.

225 Pointer equality should not be confused with alias analysis: we could require == to return true for pointers with 226 the same address but different provenance, while still permitting alias analysis to regard the two as distinct by 227 making accesses via pointers with the wrong provenance UB. 228

**Pointer relational comparison and provenance** In ISO C (6.5.8p5), inter-object pointer relational compari-229 son (with < etc.) is undefined behaviour. Just as for inter-object pointer subtraction, there are platforms where this 230 would go wrong, but there are also substantial bodies of code that rely on it, e.g. for lock orderings 231

It may be desirable to make it implementation-defined whether such pointer construction is allowed.

234 235

232 233

221

222

223

- 236
- 237
- 238 239
- 240

#### **REFINING THE BASIC PROVENANCE MODEL TO SUPPORT POINTER CONSTRUCTION VIA** 3 241 CASTS, REPRESENTATION ACCESSES, ETC. 242

To support low-level systems programming, C provides many other ways to construct and manipulate pointer 243 values: 244

- casts of pointers to integer types and back, possibly with integer arithmetic, e.g. to force alignment, or to store information in unused bits of pointers;
- copying pointer values with memcpy;

245

246

247

250

254

255

259

260

261

262

263

264

265

266

267

268

269

270

271

276

277

278

279

280

281

282 283

284

285

286

287

288

290

- 248 • manipulation of the representation bytes of pointers, e.g. via user code that copies them via char\* or 249 unsigned char\* accesses;
  - type punning between pointer and integer values;
- 251 • I/O, using either fprintf/fscanf and the %p format, fwrite/fread on the pointer representation bytes, or 252 pointer/integer casts and integer I/O; 253
  - 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.

256 A satisfactory semantics has to address all these, together with the implications on optimisation. We define and 257 explore several alternatives: 258

- PNVI-plain: a semantics that does not track provenance via integers, but instead, at integer-to-pointer cast points, checks whether the given address points within a live object and, if so, recreates the corresponding provenance. We explain in the next section why this is not as damaging to optimisation as it may sound.
- PNVI-ae (PNVI exposed-address): a variant of PNVI that allows integer-to-pointer casts to recreate provenance only for storage instances that have previously been exposed. A storage instance is deemed exposed by a cast of a pointer to it to an integer type, by a read (at non-pointer type) of the representation of the pointer, or by an output of the pointer using %p.
  - PNVI-ae-udi (PNVI exposed-address user-disambiguation): a further refinement of PNVI-ae that supports roundtrip casts, from pointer to integer and back, of pointers that are one-past a storage instance. This is the currently preferred option in the C memory object model study group.
- PVI: a semantics that tracks provenance via integer computation, associating a provenance with all integer values (not just pointer values), preserving provenance through integer/ pointer casts, and making some particular choices for the provenance results of integer and pointer +/integer operations; or

272 We write PNVI-\* for PNVI-plain, PNVI-ae, and PNVI-ae-udi. The PNVI-plain and PVI semantics were described in 273 the POPL 2019/N2311 paper. PNVI-ae and PNVI-ae-udi have emerged from discussions in the C memory object 274 model study group. 275

We also mention other variants of PNVI that seem less desirable:

- PNVI-address-taken: an earlier variant of PNVI-ae that allowed integer-to-pointer casts to recreate provenance for objects whose address has been taken (irrespective of whether it has been exposed); and
- PNVI-wildcard: a variant that gives a "wildcard" provenance to the results of integer-to-pointer casts, delaying checks to access time.

The PVI semantics, originally developed informally in ISO WG14 working papers [Memarian et al. 2018; Memarian and Sewell 2016a], was motivated in part by the GCC documentation [FSF 2018]:

"When casting from pointer to integer and back again, the resulting pointer must reference the same object as the original pointer, otherwise the behavior is undefined. That is, one may not use integer arithmetic to avoid the undefined behavior of pointer arithmetic as proscribed in C99 and C11 6.5.6/8."

which presumes there is an "original" pointer, and by experimental data for uintptr\_t analogues of the first test of §2, which suggested that GCC and ICC sometimes track provenance via integers (see xy and yx variants). However, discussions at the 2018 GNU Tools Cauldron suggest instead that at least some key developers regard the result 289 of casts from integer types as potentially broadly aliasing, at least in their GIMPLE IR, and such test results as long-standing bugs in the RTL backend.

- 291 292
- 293
- 294
- 295 296

299

- 297
- 298

305

306

307

308

309

310

311

312

313

314 315

338

### 4 REFINING THE BASIC PROVENANCE MODEL: PHENOMENA AND EXAMPLES

<sup>302</sup> <sup>303</sup> **Pointer/integer casts** The ISO standard (6.3.2.3) leaves conversions between pointer and integer types almost <sup>304</sup> entirely implementation-defined, except for conversion of integer constant 0 and null pointers, and for the

optional intptr\_t and uintptr\_t types, for which // provenance\_roundtrip\_via\_intptr\_t.c it guarantees that any "valid pointer to void" can #include <stdio.h> be converted and back, and that "the result will #include <inttypes.h> *compare equal to the original pointer*". As we have int x=1: seen, in a post-DR260 CR provenance-aware seint main() { mantics, "compare equal" is not enough to guar**int** \*p = &x;intptr\_t i = (intptr\_t)p; antee the two are interchangeable, which was **int** \*q = (**int** \*)i; clearly the intent of that phrasing. All variants of \*q = 11; // is this free of undefined behaviour? PNVI-\* and PVI support this, by reconstructing or printf("\*p=%d \*q=%d\n",\*p,\*q); preserving the original provenance respectively.

**Inter-object integer arithmetic** Below is a uintptr\_t analogue of the §2 example pointer\_offset\_from\_ptr\_subtraction\_global\_xy.c, attempting to move between objects with uintptr\_t

arithmetic. In PNVI-\*, this has defined behaviour. 318 For PNVI-plain: the integer values are pure inte-319 gers, and at the int \* cast the value of ux+offset 320 matches the address of y (live and of the right 321 type), so the resulting pointer value takes on the 322 provenance of the y storage instance. For PNVI-323 ae and PNVI-ae-udi, the storage instance for y is 324 marked as *exposed* at the cast of &y to an integer, 325 and so the above is likewise permitted there. 326

In PVI, this is UB. First, the integer values of 327 328 ux and uy have the provenances of the storage instances of x and y respectively. Then offset 329 is a subtraction of two integer values with non-330 equal single provenances; we define the result 331 332 of such to have the empty provenance. Adding that empty-provenance result to ux preserves the 333 original x-storage instance provenance of the lat-334 ter, as does the cast to int\*. Then the final \*p=11 } 335

```
// pointer_offset_from_int_subtraction_global_xy.c
#include <stdio.h>
#include <string.h>
#include <stdint.h>
#include <inttypes.h>
int x=1, y=2;
int main() {
  uintptr_t ux = (uintptr_t)&x;
  uintptr_t uy = (uintptr_t)&y;
  uintptr_t offset = uy - ux;
  printf("Addresses: &x=%"PRIuPTR" &y=%"PRIuPTR\
         " offset=%"PRIuPTR" \n",ux,uy,offset);
  int *p = (int *)(ux + offset);
  int *q = &y;
  if (memcmp(&p, &q, sizeof(p)) == 0) {
    *p = 11; // is this free of UB?
    printf("x=%d y=%d *p=%d *q=%d\n",x,y,*p,*q);
  }
```

access is via a pointer value whose address is not consistent with its provenance. Similarly, PNVI-\* allows (contrary
 to current GCC/ICC O2) a uintptr\_t analogue of the first test of §2, on the left below. PVI forbids this test.

```
339
       // provenance_basic_using_uintptr_t_global_xy.c
                                                            // pointer_offset_xor_global.c
340
                                                            #include <stdio.h>
       #include <stdio.h>
341
       #include <string.h>
                                                            #include <inttypes.h>
342
       #include <stdint.h>
                                                            int x=1;
       #include <inttypes.h>
                                                            int y=2;
343
       int x=1, y=2;
                                                            int main() {
344
                                                              int *p = &x;
       int main() {
345
         uintptr_t ux = (uintptr_t)&x;
                                                              int *q = &y;
346
         uintptr_t uy = (uintptr_t)&y;
                                                              uintptr_t i = (uintptr_t) p;
347
         uintptr_t offset = 4;
                                                               uintptr_t j = (uintptr_t) q;
348
         ux = ux + offset;
                                                               uintptr_t k = i ^ j;
349
         int *p = (int *)ux; // does this have UB?
                                                               uintptr_t l = k ^ i;
350
                                                               int *r = (int *)l;
         int *q = &y;
351
         printf("Addresses: &x=%p p=%p &y=%"PRIxPTR\
                                                               // are r and q now equivalent?
                 "\n",(void*)&x,(void*)p,uy);
                                                               *r = 11;
                                                                            // does this have defined behaviour?
352
         if (memcmp(&p, &q, sizeof(p)) == 0) {
                                                               _Bool b = (r==q);
353
           *p = 11; // does this have undefined behaviour?
                                                              printf("x=%i y=%i *r=%i (r==p)=%s\n",x,y,*r,
354
           printf("x=%d y=%d *p=%d
                                                                      b?"true":"false");
355
       *q=%d\n",x,y,*p,*q);
                                                            }
356
        }
357
       }
358
359
```

```
360
```

Both choices are defensible here: PVI will permit more aggressive alias analysis for pointers computed via integers (though those may be relatively uncommon), while PNVI-\* will allow not just this test, which as written is probably not idiomatic desirable C, but also the essentially identical XOR doubly linked list idiom, using only one pointer per node by storing the XOR of two, on the right above. Opinions differ as to whether that idiom matters for modern code.

There are other real-world but rare cases of inter-object arithmetic, e.g. in the implementations of Linux and FreeBSD per-CPU variables, in fixing up pointers after a realloc, and in dynamic linking (though arguably some of these are not between C abstract-machine objects). These are rare enough that it seems reasonable to require additional source annotation, or some other mechanism, to prevent compilers implicitly assuming that uses of such pointers as undefined.

Pointer provenance for pointer bit manipulations It is a standard idiom in systems code
 to use otherwise unused bits of pointers: low-order bits for pointers known to be aligned,
 and/ar high order bits havend the addressed la

```
and/or high-order bits beyond the addressable
                                                     // provenance_tag_bits_via_uintptr_t_1.c
374
      range. The example on the right (which as-
                                                     #include <stdio.h>
375
      sumes _Alignof(int) >= 4) does this: casting
                                                     #include <stdint.h>
376
      a pointer to uintptr_t and back, using bitwise
                                                     int x=1;
377
      logical operations on the integer value to store
                                                     int main() {
378
      some tag bits.
                                                       int *p = \&x;
379
                                                       // cast &x to an integer
        To allow this, we suggest that the set of un-
                                                       uintptr_t i = (uintptr_t) p;
380
      used bits for pointer types of each alignment
                                                       // set low-order bit
381
      should be made implementation-defined. In
                                                       i = i | 1u;
382
      PNVI-* the intermediate value of q will have
                                                       // cast back to a pointer
383
      empty provenance, but the value of r used for
                                                       int *q = (int *) i; // does this have UB?
384
      the access will re-acquire the correct prove-
                                                       // cast to integer and mask out low-order bits
385
      nance at cast time. In PVI we make the binary
                                                       uintptr_t j = ((uintptr_t)q) & ~((uintptr_t)3u);
386
      operations used here, combining an integer
                                                       // cast back to a pointer
387
      value that has some provenance ID with a pure
                                                       int *r = (int *) j;
388
      integer, preserve that provenance.
                                                       // are r and p now equivalent?
389
                                                                           // does this have UB?
                                                       *r = 11;
        (A separate question is the behaviour if
                                                       _Bool b = (r==p); // is this true?
390
      the integer value with tag bits set is con-
                                                       printf("x=%i *r=%i (r==p)=%s\n",x,*r,b?"t":"f");
391
      verted back to pointer type. In ISO the result is
                                                     }
392
      implementation-defined, per 6.3.2.3p{5,6} and
```

<sup>393</sup> 7.20.1.4.)

394

Algebraic properties of integer operations The PVI definitions of the provenance results of integer operations, chosen to make pointer\_offset\_from\_int\_subtraction\_global\_xy.c forbidden and provenance\_tag\_bits\_via\_uintptr\_t\_1.c allowed, has an unfortunate consequence: it makes those operations no longer associative. Compare the examples below:

```
399
                                   // pointer_arith_algebraic_properties_2_global.c
400
      #include <stdio.h>
401
      #include <inttypes.h>
402
      int y[2], x[2];
      int main() {
403
        int *p=(int*)(((uintptr_t)&(x[0])) +
404
          (((uintptr_t)&(y[1]))-((uintptr_t)&(y[0]))));
405
        *p = 11; // is this free of undefined behaviour?
406
        printf("x[1]=%d *p=%d\n",x[1],*p);
407
        return 0;
408
      }
409
                                   // pointer_arith_algebraic_properties_3_global.c
410
      #include <stdio.h>
411
      #include <inttypes.h>
412
      int y[2], x[2];
413
      int main() {
        int *p=(int*)(
414
          (((uintptr_t)&(x[0])) + ((uintptr_t)&(y[1])))
415
          -((uintptr_t)&(y[0])) );
416
        *p = 11; // is this free of undefined behaviour?
417
        //(equivalent to the &x[0]+(&(y[1])-&(y[0])) version?)
418
        printf("x[1]=%d *p=%d\n",x[1],*p);
419
420
```

}

#### return 0;

421 422 423

428

441

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

come permitted for all storage instances for

The latter is UB in PVI. It is unclear whether this would be acceptable in practice, either for C programmers or for compiler optimisation. One could conceivably switch to a PVI-multiple variant, allowing provenances to be finite sets of storage instance IDs. That would allow the pointer\_offset\_from\_int\_subtraction\_global\_xy.c example above, but perhaps too much else besides. The PNVI-\* models do not suffer from this problem.

429
 429
 429
 430
 430
 430
 430
 430
 430
 430
 431
 431
 432
 432
 433
 434
 434
 435
 435
 436
 437
 438
 438
 439
 430
 430
 430
 430
 430
 431
 431
 432
 432
 433
 434
 435
 435
 436
 437
 438
 438
 438
 438
 438
 439
 430
 430
 430
 430
 430
 431
 431
 432
 431
 432
 432
 432
 433
 434
 434
 435
 435
 435
 435
 435
 435
 436
 437
 437
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 438
 439
 439
 439
 439

```
The ISO C18 text does not explicitly address this
431
                                                         // pointer_copy_memcpy.c
      (in a pre-provenance semantics, before DR260,
432
                                                        #include <stdio.h>
      it did not need to). One could do so by special-
                                                        #include <string.h>
433
      casing memcpy() and similar functions to preserve
                                                         int x=1:
434
      provenance, but the following questions suggest
                                                         int main() {
435
      less ad hoc approaches, for PNVI-plain or PVI.
                                                           int *p = &x;
436
                                                           int *q;
      For PNVI-ae and PNVI-ae-udi, the best approach
437
                                                           memcpy (&q, &p, sizeof p);
      is not yet clear.
438
                                                           *q = 11; // is this free of undefined behaviour?
439
                                                           printf("*p=%d *q=%d\n",*p,*q);
440
                                                        }
```

442 Copying pointer values bytewise, with user-memcpy
 443 Unit of object representations, e.g. as in the following naive user implementation of a memcpy-like function,

```
which constructs a pointer value from copied
                                              // pointer_copy_user_dataflow_direct_bytewise.c
bytes. This too should be allowed. PNVI-plain
                                              #include <stdio.h>
makes it legal: the representation bytes have
                                              #include <string.h>
no provenance, but when reading a pointer
                                              int x=1;
value from the copied memory, the read will
                                              void user_memcpy(unsigned char* dest,
be from multiple representation-byte writes.
                                                                unsigned char *src, size_t n) {
                                                while (n > 0)
We use essentially the same semantics for such
                                                                {
                                                   *dest = *src;
reads as for integer-to-pointer casts: checking
                                                   src += 1; dest += 1; n -= 1;
at read-time that the address is within a live
                                                }
object, and giving the result the correspond-
ing provenance. For PNVI-ae and PNVI-ae-udi,
                                              int main() {
the current proposal is to mark storage in-
                                                int *p = &x;
stances as exposed whenever representation
                                                int *q;
bytes of pointers to them are read, and use
                                                user_memcpy((unsigned char*)&q,
                                                             (unsigned char*)&p, sizeof(int *));
the same semantics for reads of pointer values
                                                 *q = 11; // is this free of undefined behaviour?
from representation-byte writes as for integer-
                                                printf("*p=%d *q=%d\n",*p,*q);
to-pointer casts. This is attractively simple, but
it does means that integer-to-pointer casts be-
```

which a pointer has been copied via user\_memcpy, which is arguably too liberal. It may be possible to add additional annotations for code like user\_memcpy to indicate (to alias analysis) that (a) their target memory should have the same provenance as their source memory, and (b) the storage instances of any copied pointers should not be marked as exposed, despite the reads of their representation bytes. This machinery has not yet been designed.

467 One might instead think of recording symbolically in the semantics of integer values (e.g. for representation-byte 468 values) whether they are of the form "byte *n* of pointer value *v*", or perhaps "byte *n* of pointer value of type *t*", 469 and allow reads of pointer values from representation-byte writes only for such. This is more complex and rather 470 ad hoc, arbitrarily restricting the integer computation that can be done on such bytes. If one wanted to allow (e.g.) 471 bitwise operations on such bytes, as in provenance\_tag\_bits\_via\_repr\_byte\_1.c, one would essentially have to 472 adopt a PVI model. However, note that to capture the 6.5p6 preservation of effective types by character-type array 473 copy ("If a value is copied into an object having no declared type using memcpy or memmove, or is copied as an array 474 of character type, then the effective type of the modified object for that access and for subsequent accesses that do not 475 modify the value is the effective type of the object from which the value is copied, if it has one."), we might need 476 something like a very restricted version of PVI: some effective-type information attached to integer values of 477 character type, to say "byte n of pointer value of type t", with all integer operations except character-type stores 478 clearing that info. 479

#### C provenance semantics: examples

As Lee observes [private communication], to make it legal for compilers to replace user-memcpy by the library version, one might want the two to have exactly the same semantics. Though strictly speaking that is a question about the compiler intermediate language semantics, not C source semantics.

PVI makes user-memcpy legal by regarding each byte (as an integer value) as having the provenance of the original pointer, and the result pointer, being composed of representation bytes of which at least one has that provenance and none have a conflicting provenance, as having the same.

Real memcpy() implementations are more complex. The glibc memcpy()[glibc 2018] involves copying byte-by-byte, as above, and also word-by-word and, using virtual memory manipulation, page-by-page. Word-by-word copying is not permitted by the ISO standard, as it violates the effective type rules, but we believe C2x should support it for suitably annotated code. Virtual memory manipulation is outside our scope at present.

Reading pointer values from byte writes In all these provenance semantics, pointer values carry their provenance unchanged, both while manipulated in expressions (e.g. with pointer arithmetic) and when stored or loaded as values of pointer type. In the detailed semantics, memory contains abstract bytes rather than general C language values, and so we record provenance in memory by attaching a provenance to each abstract byte. For pointer values stored by single writes, this will usually be identical in each abstract byte of the value.

496 However, we also have to define the result of reading a pointer value that has been partially or completely 497 written by (integer) representation-byte writes. In PNVI-\*, we use the same semantics as for integer-to-pointer 498 casts, reading the numeric address and reconstructing the associated provenance iff a live storage instance covering 499 that address exists (and, for PNVI-ae and PNVI-ae-udi, if that instance has been exposed). To determine whether 500 a pointer value read is from a single pointer value write (and thus should retain its original provenance when 501 read), or from a combination of representation byte writes and perhaps also a pointer value write (and thus 502 should use the integer-to-pointer cast semantics when read), we also record, in each abstract byte, an optional 503 pointer-byte index (e.g. in 0..7 on an implementation with 8-byte pointer values). Pointer value writes will set 504 these to the consecutive sequence 0, 1, ..., 7, while other writes will clear them. For example, the code on the left 505 below sets the fourth byte of p to 0. The memory state on the right, just after the \*q=2, shows the pointer-byte 506 indices of p, one of which has been cleared (shown as -). When the value of p is read (e.g. in the q=p), the fact that 507 there is not a consecutive sequence 0, 1, ..., 7 means that PNVI-\* will apply the integer-to-pointer cast semantics, 508 here successfully recovering the provenance @68 of the storage instance x. Then the write of q will itself have a 509 consecutive sequence (its pointer-byte indices are therefore suppressed in the diagram). Any non-pointer write 510 overlapping the footprint of p, or any pointer write that overlaps that footprint but does not cover it all, would 511 interrupt the consecutive sequence of indices. 512



In PNVI-plain a representation-byte copy of a pointer value thus is subtly different from a copy done at pointer
 type: the latter retains the original provenance, while the former, when it is loaded, will take on the provenance of
 whatever storage instance is live (and covers its address) *at load time*.

540

The conditional in the example is needed to avoid UB: the semantics does not constrain the allocation address of x, so there are executions in which byte 4 is not 0, in which case the read of p would have a wild address and the empty provenance, and the write \*q=2 would flag UB.

Pointer provenance for bytewise pointer representation manipulations To examine the possible seman tics for pointer representation bytes more closely, especially for PNVI-ae and PNVI-ae-udi, consider the following.
 As in provenance\_tag\_bits\_via\_uintptr\_t\_1.c, it manipulates the low-order bits of a pointer value, but now it
 does so by manipulating one of its representation bytes (as in pointer\_copy\_user\_dataflow\_direct\_bytewise.c)

```
548
      instead of by casting to uintptr_t
                                          // provenance_tag_bits_via_repr_byte_1.c
549
      and back. In PNVI-plain and PVI
                                         #include <assert.h>
550
      this will just work, respectively
                                         #include <stdio.h>
551
      reconstructing the original prove- #include <stdint.h>
552
      nance and tracking it through the int x=1;
553
      (changed and unchanged) integer
                                         int main() {
554
                                            int *p=&x, *q=&x;
      bytes.
                                            // read low-order (little endian) representation byte of p
555
        In PNVI-ae and PNVI-ae-udi, we
                                            unsigned char i = *(unsigned char*)&p;
556
      regard the storage instance of x as
                                            // check the bottom two bits of an int* are not used
557
      having been exposed by the read
                                            assert(_Alignof(int) >= 4);
558
      of a pointer value (with non-empty
                                            assert((i & 3u) == 0u);
559
      provenance in its abstract bytes in
                                            // set the low-order bit of the byte
560
      memory) at an integer (really, non-
                                            i = i | 1u;
561
      pointer) type. Then the last reads
                                            // write the representation byte back
562
      of the value of p, from a combi-
                                            *(unsigned char*)&p = i;
563
                                            // [p might be passed around or copied here]
      nation of the original p=&x write
564
                                            // clear the low-order bits again
      and later integer byte writes, use
                                            *(unsigned char*)\&p = (*(unsigned char*)\&p) \& ~((unsigned char)3u);
565
      the same semantics as integer-to-
                                            // are p and q now equivalent?
566
      pointer casts, and thus recreate the
                                            *p = 11;
                                                              // does this have defined behaviour?
567
      original provenance.
                                            _Bool b = (p==q); // is this true?
568
                                            printf("x=%i *p=%i (p==q)=%s\n",x,*p,b?"true":"false");
569
                                          }
570
571
```

**Copying pointer values via encryption** To more clearly delimit what idioms our proposals do and do not allow, consider copying pointers via code that encrypts or compresses a block of multiple pointers together, decrypting or uncompressing later.

In PNVI-plain, it would just work, in the same way as user\_memcpy(). In PNVI-ae and PNVI-ae-udi, it would work but leave storage instances pointed to by those pointers exposed (irrespective of whether the encryption is done via casts to integer types or by reads of representation bytes), similar to user\_memcpy and provenance\_tag\_bits\_via\_repr\_byte\_1.c.

One might argue that pointer construction via intptr\_t and back via any value-dependent identity function should be required to work. That would admit these, but defining that notion of "value-dependent" is exactly what is hard in the concurrency thin-air problem [Batty et al. 2015], and we do not believe that it is practical to make compilers respect dependencies in general.

In PVI, this case involves exactly the same combination of distinct-provenance values that (to prohibit interobject arithmetic, and thereby enable alias analysis) we above regard as having empty-provenance results. As copying pointers in this way is a very rare idiom, one can argue that it is reasonable to require such code to have additional annotations.

593

578

579

580

581

582

583

- 594
- 595
- 596
- 597
- 598
- 599 600

637

638

639

660

601 **Copying pointer values via control flow** We also have to ask whether a usable pointer can be constructed 602 via non-dataflow control-flow paths, e.g. if testing equality of an unprovenanced integer value against a valid 603 pointer permits the integer to be used as if it had the same provenance as the pointer. We do not believe that this 604 is relied on in practice. For example, consider exotic versions of memcpy that make a control-flow choice on the 605 value of each bit or each byte, reconstructing each with constants in each control-flow branch

```
// pointer_copy_user_ctrlflow_bytewise_abbrev.c
                                                             // pointer_copy_user_ctrlflow_bitwise.c
607
       #include <stdio.h>
                                                             #include <stdio.h>
608
       #include <string.h>
                                                             #include <inttypes.h>
609
       #include <assert.h>
                                                             #include <limits.h>
610
       #include <limits.h>
                                                             int x=1;
611
                                                             int main() {
       int x=1;
612
       unsigned char control_flow_copy(unsigned char c) {
                                                               int *p = \&x;
613
         assert(UCHAR_MAX==255);
                                                               uintptr_t i = (uintptr_t)p;
         switch (c) {
                                                               int uintptr_t_width = sizeof(uintptr_t) * CHAR_BIT;
614
         case 0: return(0);
                                                               uintptr_t bit, j;
615
         case 1: return(1);
                                                               int k;
616
         case 2: return(2);
                                                               i=0;
617
                                                                for (k=0; k<uintptr_t_width; k++) {</pre>
         . . .
618
         case 255: return(255);
                                                                  bit = (i & (((uintptr_t)1) << k)) >> k;
619
         }
                                                                  if (bit == 1)
620
       }
                                                                    j = j | ((uintptr_t)1 << k);</pre>
621
       void user_memcpy2(unsigned char* dest,
                                                                  else
622
                          unsigned char *src, size_t n) {
                                                                    j = j;
623
         while (n > 0) {
                                                               }
624
           *dest = control_flow_copy(*src);
                                                               int *q = (int *)j;
                                                                *q = 11; // is this free of undefined behaviour?
625
           src += 1;
           dest += 1;
                                                               printf("*p=%d *q=%d\n",*p,*q);
626
           n -= 1;
                                                             }
627
         }
628
       }
629
       int main() {
630
         int *p = &x;
631
         int *q;
632
         user_memcpy2((unsigned char*)&q, (unsigned char*)&p,
633
                      sizeof(p));
634
         *q = 11; // does this have undefined behaviour?
635
         printf("*p=%d *q=%d\n",*p,*q);
       }
636
```

In PNVI-plain these would both work. In PNVI-ae and PNVI-ae-udi they would also work, as the first exposes the storage instance of the copied pointer value by representation-byte reads and the second by a pointer-to-integer cast. In PVI they would give empty-provenance pointer values and hence UB.

Integer comparison and provenance If integer values have associated provenance, as in PVI, one has to ask whether the result of an integer comparison should also be allowed to be provenance dependent (provenance\_equality\_uintptr\_t\_global\_xy.c). GCC did do so at one point, but it was regarded as a bug and fixed (from 4.7.1 to 4.8). We propose that the numeric results of all operations on integers should be unaffected by the provenances of their arguments. For PNVI-\*, this question is moot, as there integer values have no provenance.

Pointer provenance and union type punning Pointer values can also be constructed in C by type punning, e.g. writing a pointer-type union member, reading it as a uintptr\_t union member, and then casting back to a pointer type. (The example assumes that the object representation of the pointer and the object representation of the result of the cast to integer are identical. This property is not guaranteed by the C standard, but holds for many implementations.)

The ISO standard says "the appropriate part of the object representation of the value is reinterpreted as an object representation in the new type", but says little about that reinterpretation. We propose that these reinterpretations be required to be implementation-defined, and, in PNVI-plain, that the usual integer-to-pointer cast semantics be used at such reads.

For PNVI-ae and PNVI-ae-udi, the same semantics as for
representation-byte reads also permits this case: the storage
instance is deemed to be 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.
For PVI we propose that it be implementation-defined

For PVI, we propose that it be implementation-defined whether the result preserves the original provenance } (e.g. where they are the identity).

```
// provenance_union_punning_3_global.c
#include <stdio.h>
#include <stdio.h>
#include <string.h>
#include <inttypes.h>
int x=1;
typedef union { uintptr_t ui; int *up; } un;
int main() {
    un u;
    int *p = &x;
    u.up = p;
    uintptr_t i = u.ui;
    int *q = (int*)i;
    *q = 11; // does this have UB?
    printf("x=%d *p=%d *q=%d\n",x,*p,*q);
    return 0;
}
```

**Pointer provenance via IO** Consider now pointer provenance flowing via IO, e.g. writing the address of an object to a string, pipe or file and reading it back in. We have three versions: one using fprintf/fscanf and the %p format, one using fwrite/fread on the pointer representation bytes, and one converting the pointer to and from uintptr\_t and using fprintf/fscanf on that value with the PRIuPTR/SCNuPTR formats (provenance\_via\_io\_percentp\_global.c, provenance\_via\_io\_bytewise\_global.c, and provenance\_via\_io\_uintptr\_t\_global.c) The first gives a syntactic indication of a potentially escaping pointer value, while the others (after preprocessing) do not. Somewhat exotic though they are, these idioms are used in practice: in graphics code for serialisation/deserialisation (using %p), in xlib (using SCNuPTR), and in debuggers.

In the ISO standard, the text for fprintf and scanf for %p says that this should work: "If the input item is a value converted earlier during the same program execution, the pointer that results shall compare equal to that value; otherwise the behavior of the %p conversion is undefined" (again construing the pre-DR260 "compare equal" as implying the result should be usable for access), and the text for uintptr\_t and the presence of SCNuPTR in inttypes.h weakly implies the same there.

But then what can compiler alias analyses assume about such a pointer read? In PNVI-plain, this is simple: at scanf-time, for the %p version, or when a pointer is read from memory written by the other two, we can do a runtime check and potential acquisition of provenance exactly like an integer-to-pointer cast.

In PNVI-ae and PNVI-ae-udi, for the %p case we mark the associated storage instance as exposed by the output, and use the same semantics as integer-to-pointer casts on the input. The uintptr\_t case and representation-byte case also mark the storage instance as exposed, in the normal way for these models.

For PVI, there are several options, none of which seem ideal: we could use a PNVI-like semantics, but that would be stylistically inconsistent with the rest of PVI; or (only for the first) we could restrict that to provenances

717

718

685

686

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

that have been output via %p), or we could require new programmer annotation, at output and/or input points, to 721 constrain alias analysis. 722

723 **Pointers from device memory and linking** In practice, concrete memory addresses or relationships between 724 them sometimes are determined and relied on by programmers, in implementation-specific ways. Sometimes 725 these are simply concrete absolute addresses which will never alias C stack, heap, or program memory, e.g. those 726 of particular memory-mapped devices in an embedded system. Others are absolute addresses and relative layout 727 of program code and data, usually involving one or more linking steps. For example, platforms may lay out certain 728 regions of memory so as to obey particular relationships, e.g. in a commodity operating system where high 729 addresses are used for kernel mappings, initial stack lives immediately below the arguments passed from the 730 operating system, and so on. The details of linking and of platform memory maps are outside the scope of ISO C, 731 but real C code may embody knowledge of them. Such code might be as simple as casting a platform-specified 732 address, represented as an integer literal, to a pointer. It might be more subtle, such as assuming that one object 733 directly follows another in memory—the programmer having established this property at link time (perhaps by a 734 custom linker script). It is necessary to preserve the legitimacy of such C code, so that compilers may not view 735 such memory accesses as undefined behaviour, even with increasing link-time optimisation.

736 We leave the design of exactly what escape-hatch mechanisms are needed here as an open problem. For 737 memory-mapped devices, one could simply posit implementation-defined ranges of such memory which are 738 guaranteed not to alias C objects. The more general linkage case is more interesting, but well outside current ISO 739 C. The tracking of provenance through embedded assembly is similar. 740

Pointers from allocator libraries Our semantics special-cases malloc and the related functions, by giving 741 their results fresh provenances. This is stylistically consistent with the ISO text, which also special-cases them, but 742 it would be better for C to support a general-purpose annotation, to let both stdlib implementations and other 743 libraries return pointers that are treated as having fresh provenance outside (but not inside) their abstraction 744 boundaries. 745

Compilers already have related annotations, e.g. GCC's malloc attribute "tells the compiler that a function is 746 malloc-like, i.e., that the pointer P returned by the function cannot alias any other pointer valid when the function 747 returns, and moreover no pointers to valid objects occur in any storage addressed by P" (https://gcc.gnu.org/onlinedocs/ 748 gcc/Common-Function-Attributes.html#Common-Function-Attributes). 749

750 751

780

#### IMPLICATIONS OF PROVENANCE SEMANTICS FOR PPTIMISATIONS 5

752 In an ideal world, a memory object semantics for C would be consistent with all existing mainstream code usage 753 and compiler behaviour. In practice, we suspect that (absent a precise standard) these have diverged too much 754 for that, making some compromise required. As we have already seen, the PNVI semantics would make some 755 currently observed GCC and ICC behaviour unsound, though at least some key GCC developers already regard 756 that behaviour as a longstanding unfixed bug, due to the lack of integer/pointer type distinctions in RTL. We now 757 consider some other important cases, by example. 758

**Optimisation based on equality tests** Both PNVI-\* and PVI let p==q hold in some cases where p and q are 759 not interchangeable. As Lee et al. [2018] observe in the LLVM IR context, that may limit optimisations such as 760 GVN (global value numbering) based on pointer equality tests. PVI suffers from the same problem also for integer 761 comparisons, wherever the integers might have been cast from pointers and eventually be cast back. This may be 762 more serious. 763

```
764
      Can a function argument alias local variables of the function?
765
      this to be forbidden, to let optimisation assume its absence. Consider first the example below, where
766
      main() guesses the address of f()'s local variable, passing it in as
767
      a pointer, and f() checks it before using it for an access. Here we
768
      see, for example, GCC -O0 optimising away the if and the write
769
      *p=7, even in executions where the ADDRESS_PFI_1PG constant is
770
      the same as the printf'd address of j. We believe that compiler
771
      behaviour should be permitted, and hence that this program should
772
      be deemed to have UB - or, in other words, that code should not
773
      normally be allowed to rely on implementation facts about the
774
      allocation addresses of C variables.
```

775 The PNVI-\* semantics deems this to be UB, because at the point 776 of the (int\*)i cast the j storage instance does not vet exist (let 777 alone, for PNVI-ae and PNVI-ae-udi, having been exposed by 778 having one of its addresses taken and cast to integer), so the cast 779

In general one would like // pointer\_from\_integer\_lpg.c

```
#include <stdio.h>
#include <stdint.h>
#include "charon_address_guesses.h"
void f(int *p) {
  int j=5;
  if (p==&j)
    *p=7;
  printf("j=%d &j=%p\n",j,(void*)&j);
}
int main() {
  uintptr_t i = ADDRESS_PFI_1PG;
  int *p = (int*)i;
  f(p);
}
                         Draft of April 1, 2019
```

gives a pointer with empty provenance; any execution that goes 781 into the if would thus flag UB. The PVI semantics flags UB for the 782 simple reason that j is created with the empty provenance, and 783 hence p inherits that. 784

Varying to do the cast to int \* in f() instead of main(), passing 785 in an integer i instead of a pointer, this becomes defined in PNVI-786 plain, as j exists at the point when the abstract machine does 787 the (int\*)i cast. But in PNVI-ae and PNVI-ae-udi, the storage 788 instance of j is not exposed, so the cast to int\* gives a pointer 789 with empty provenance and the access via it is UB. This example 790 is also UB in PVI. 791

At present we do not see any strong reason why making this de-792 fined would not be acceptable – it amounts to requiring compilers 793 to be conservative for the results of integer-to-pointer casts where 794 they cannot see the source of the integer, which we imagine to be 795 a rare case - but this does not match current O2 or O3 compilation 796 for GCC, Clang, or ICC. 797

```
// pointer_from_integer_lig.c
#include <stdio.h>
#include <stdint.h>
#include "charon_address_guesses.h"
void f(uintptr_t i) {
  int j=5;
  int *p = (int*)i;
  if (p==&j)
    *p=7:
  printf("j=%d &j=%p\n",j,(void*)&j);
int main() {
  uintptr_t j = ADDRESS_PFI_1IG;
  f(j);
}
```

Allocation-address nondeterminism Note that both of the previous examples take the address of j to guard their \*p=7 accesses. Removing the conditional guards gives the left and middle tests below, that one would surely like to forbid:

803			
804	<pre>// pointer_from_integer_1p.c</pre>	<pre>// pointer_from_integer_li.c</pre>	<pre>// pointer_from_integer_lie.c</pre>
805	<pre>#include <stdio.h></stdio.h></pre>	<pre>#include <stdio.h></stdio.h></pre>	<pre>#include <stdio.h></stdio.h></pre>
805	<pre>#include <stdint.h></stdint.h></pre>	<pre>#include <stdint.h></stdint.h></pre>	<pre>#include <stdint.h></stdint.h></pre>
806	<pre>#include "charon_address_guesses</pre>	.h' <b>#include</b> "charon_address_guesses	.h' <b>#include</b> "charon_address_guesses.h'
807	<pre>void f(int *p) {</pre>	<pre>void f(uintptr_t i) {</pre>	<pre>void f(uintptr_t i) {</pre>
808	<b>int</b> j=5;	<b>int</b> j=5;	<b>int</b> j=5;
809	*p=7;	<pre>int *p = (int*)i;</pre>	uintptr_t k = (uintptr_t)&j
810	printf("j=%d\n",j);	*p=7;	<pre>int *p = (int*)i;</pre>
811	}	printf("j=%d\n",j);	*p=7;
812	<pre>int main() {</pre>	}	printf("j=%d\n",j);
813	uintptr_t i = ADDRESS_PFI_1P;	<pre>int main() {</pre>	}
814	<b>int</b> *p = ( <b>int</b> *)i;	uintptr_t j = ADDRESS_PFI_1I;	<pre>int main() {</pre>
017	f(p);	f(j);	uintptr_t j = ADDRESS_PFI_1I;
010	}	}	f(j);
816			}

Both are forbidden in PVI for the same reason as before, and the first is forbidden in PNVI-\*, again because j does not exist at the cast point.

But the second forces us to think about how much allocation-address nondeterminism should be quantified over in the basic definition of undefined behaviour. For evaluation-order and concurrency nondeterminism, one would normally say that if there exists any execution that flags UB, then the program as a whole has UB (for the moment ignoring UB that occurs only on some paths following I/O input, which is another important question that the current ISO text does not address).

825 This view of UB seems to be unfortunate but inescapable. If one looks just at a single execution, then (at least 826 between input points) we cannot temporally bound the effects of an UB, because compilers can and do re-order 827 code w.r.t. the C abstract machine's sequencing of computation. In other words, UB may be flagged at some specific 828 point in an abstract-machine trace, but its consequences on the observed implementation behaviour might happen 829 much earlier (in practice, perhaps not very much earlier, but we do not have any good way of bounding how 830 much). But then if one execution might have UB, and hence exhibit (in an implementation) arbitrary observable 831 behaviour, then anything the standard might say about any other execution is irrelevant, because it can always be 832 masked by that arbitrary observable behaviour.

833 Accordingly, our semantics nondeterministically chooses an arbitrary address for each storage instance, subject 834 only to alignment and no-overlap constraints (ultimately one would also need to build in constraints from 835 programmer linking commands). This is equivalent to noting that the ISO standard does not constrain how 836 implementations choose storage instance addresses in any way (subject to alignment and no-overlap), and hence 837 that programmers of standard-conforming code cannot assume anything about those choices. Then in PNVI-plain, 838 the ...\_li.c example is UB because, even though there is one execution in which the guess is correct, there is 839

798 799 800

801

802

819

820

821

822

823

824

another (in fact many others) in which it is not. In those, the cast to int\* gives a pointer with empty provenance, 841 so the access flags UB – hence the whole program is UB, as desired. In PNVI-ae and PNVI-ae-udi, the ...\_1i.c 842 example is UB for a different reason: the storage instance of j is not exposed before the cast (int\*)i, and so the 843 result of that cast has empty provenance and the access \*p=7 flags UB, in every execution. However, if j is exposed, 844 as in the example on the right, these models still make it UB, now for the same reason as PNVI-plain. 845

846 **Can a function access local variables of its parent?** This too should be forbidden in general. The example 847 on the left below is forbidden by PVI, again for the simple reason that p has the empty provenance, and by 848

```
849
          // pointer_from_integer_2.c
                                                                 // pointer_from_integer_2g.c
          #include <stdio.h>
                                                                 #include <stdio.h>
850
          #include <stdint.h>
                                                                  #include <stdint.h>
851
          #include "charon_address_guesses.h"
                                                                  #include "charon_address_guesses.h"
852
          void f() {
                                                                  void f() {
853
             uintptr_t i=ADDRESS_PFI_2;
                                                                    uintptr_t i=ADDRESS_PFI_2G;
854
             int *p = (int*)i;
                                                                    int *p = (int*)i;
855
             *p=7;
                                                                    *p=7;
856
          }
                                                                  }
857
          int main() {
                                                                  int main() {
858
             int j=5;
                                                                    int i=5:
                                                                    if ((uintptr_t)&j == ADDRESS_PFI_2G)
859
             f();
             printf("j=%d\n",j);
                                                                      f();
                                                                    printf("j=%d &j=%p\n",j,(void*)&j);
          }
861
                                                                  }
862
```

PNVI-plain by allocation-address nondeterminism, as there exist abstract-machine executions in which the guessed 863 address is wrong. One cannot guard the access within f(), as the address of j is not available there. Guarding 864 the call to f() with **if** ((uintptr\_t)&j == ADDRESS\_PFI\_2) (pointer\_from\_integer\_2g.c on the right above) again 865 makes the example well-defined in PNVI-plain, as the address is correct and j exists at the int\* cast point, but 866 notice again that the guard necessarily involves &j. This does not match current Clang at O2 or O3, which print 867 j=5. 868

In PNVI-ae and PNVI-ae-udi, pointer\_from\_integer\_2.c is forbidden simply because j is never exposed 869 (and if it were, it would be forbidden for the same reason as in PNVI-plain). PNVI-ae and PNVI-ae-udi allow 870 pointer\_from\_integer\_2g.c, because the j storage instance is exposed by the (uinptr\_t)&j cast. 871

872 The PNVI-address-taken and PNVI-wildcard alternatives A different obvious refinement to PNVI would 873 be to restrict integer-to-pointer casts to recover the provenance only of objects that have had their address taken, 874 recording that in the memory state. PNVI-address-exposed is based on PNVI-address-taken but with the tighter 875 condition that the address must also have been cast to integer. 876

A rather different model is to make the results of integer-to-pointer casts have a "wildcard" provenance, deferring the check that the address matches a live object from cast-time to access-time. This would make 878 pointer\_from\_integer\_1pg.c defined, which is surely not desirable. 879

Perhaps surprisingly, the PNVI-ae and PNVI-ae-udi variants seem not to make much difference to the allowed tests, because the tests one might write tend to already be UB due to allocation-address nondeterminism, or to already take the address of an object to use it in a guard. These variants do have the conceptual advantage of identifying these UBs without requiring examination of multiple executions, but the disadvantage that whether an address has been taken is a fragile syntactic property, e.g. not preserved by dead code elimination.

884 885

877

880

881

882

- 886
- 887 888
- 880
- 800 891
- 892
- 893
- 894
- 895
- 896
- 897 898
- 899
- 900

960

Draft of April 1, 2019

The problem with lost address-takens and escapes Our PVI proposal allows computations that erase the numeric value (and hence a concrete view of the "semantic dependencies") of a pointer, but retain provenance. This makes examples like that below [Richard Smith, personal communication], in which the code correctly guesses a storage instance address (which has the empty provenance) and adds that to a zero-valued quantity (with the correct provenance), allowed in PVI. We emphasise that we do not think it especially desirable to allow such examples; this is just a consequence of choosing a straightforward provenance-via-integer semantics that allows the bytewise copying and the bitwise manipulation of pointers above. In other words, it is not clear how it could be forbidden simply in PVI.

908 However, in implementations 909 some algebraic optimisations may 910 be done before alias analysis, and 911 those optimisations might erase the 912 &x, replacing it and all the calcula-913 tion of i3 by 0x0 (a similar example 914 would have i3 = i1-i1). But then 915 alias analysis would be unable to 916 see that \*q could access x, and so re-917 port that it could not, and hence en-918 able subsequent optimisations that 919 are unsound w.r.t. PVI for this case. 920 The basic point is that whether a 921 variable has its address taken or es-922 caped in the source language is not 923 preserved by optimisation. A pos-924 sible solution, which would need 925 some implementation work for im-926

```
// provenance_lost_escape_1.c
#include <stdio.h>
#include <string.h>
#include <stdint.h>
#include "charon_address_guesses.h"
int x=1; // assume allocation ID @1, at ADDR_PLE_1
int main() {
  int *p = &x;
  uintptr_t i1 = (intptr_t)p;
                                          // (@1,ADDR_PLE_1)
  uintptr_t i2 = i1 & 0x0000000FFFFFFF;//
  uintptr_t i3 = i2 & 0xFFFFFFF00000000;// (@1,0x0)
  uintptr_t i4 = i3 + ADDR_PLE_1;
                                          // (@1, ADDR_PLE_1)
  int *q = (int *)i4;
  printf("Addresses: p=%p\n",(void*)p);
  if (memcmp(&i1, &i4, sizeof(i1)) == 0) {
    *q = 11; // does this have defined behaviour?
    printf("x=%d *p=%d *q=%d\n",x,*p,*q);
  }
}
```

plementations that do track provenance through integers, but perhaps acceptably so, would be to require those
 initial optimisation passes to record the address-takens involved in computations they erase, so that that could be
 passed in explicitly to alias analysis. In contrast to the difficulties of preserving dependencies to avoid thin-air
 concurrency, this does not forbid optimisations that remove dependencies; it merely requires them to describe
 what they do.

In PNVI-plain, the example is also allowed, but for a simpler reason that is not affected by such integer optimisation: the object exists at the **int**\* cast. Implementations that take a conservative view of all pointers formed from integers would automatically be sound w.r.t. this. At present ICC is not, at O2 or O3.

PNVI-ae and PNVI-ae-udi are more like PVI here: they allow the example, but only because the address of p is
both taken and cast to an integer type. If these semantics were used for alias analysis in an intermediate language
after such optimisation, this would likewise require the optimisation passes to record which addresses have been
taken and cast to integer (or otherwise exposed) in eliminated code, to be explicitly passed in to alias analysis.

Should PNVI allow one-past integer-to-pointer casts? For PNVI\*, one has to choose whether an integer that is one-past a live object (and not strictly within another) can be cast to a pointer with valid provenance, or whether this should give an empty-provenance pointer value. Lee observes that the latter may be necessary to make some optimisation sound [personal communication], and we imagine that this is not a common idiom in practice, so for PNVI-plain and PNVI-ae we follow the stricter semantics.

PNVI-ae-udi, however, is designed to permit a cast of a one-past pointer to integer and back to recover the
 original provenance, replacing the integer-to-pointer semantic check that *x* is properly within the footprint of the
 storage instance by a check that it is properly within or one-past. That makes the following example allowed in
 PNVI-ae-udi, while it is forbidden in PNVI-ae and PNVI-plain.

// provenance\_roundtrip\_via\_intptr\_t\_onepast.c

```
971
          #include <stdio.h>
972
          #include <inttypes.h>
          int x=1:
973
          int main() {
974
            int *p = &x;
975
            p=p+1;
976
            intptr_t i = (intptr_t)p;
977
            int *q = (int *)i;
978
            q=q-1;
979
            *q = 11; // is this free of undefined behaviour?
980
            printf("*p=%d *q=%d\n",*p,*q);
981
          }
```

The downside of this is that one has to handle pointer-to-integer casts for integer values that are ambiguously both one-past one storage instance and at the start of the next. The PNVI-ae-udi approach to that is to leave 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 the following two, each of which uses the result of the cast in one consistent way, well defined:

```
// pointer_from_int_disambiguation_1.c
                                                             // pointer_from_int_disambiguation_2.c
       #include <stdio.h>
                                                            #include <stdio.h>
989
       #include <string.h>
                                                            #include <string.h>
990
       #include <stdint.h>
                                                            #include <stdint.h>
991
       #include <inttypes.h>
                                                            #include <inttypes.h>
992
       int y=2, x=1;
                                                            int y=2, x=1;
993
       int main() {
                                                             int main() {
994
         int *p = &x+1;
                                                               int *p = &x+1;
995
         int *q = &y;
                                                               int *q = &y;
996
         uintptr_t i = (uintptr_t)p;
                                                               uintptr_t i = (uintptr_t)p;
         uintptr_t j = (uintptr_t)g;
                                                               uintptr_t j = (uintptr_t)g;
997
         if (memcmp(\&p, \&q, sizeof(p)) == 0) {
                                                               if (memcmp(\&p, \&q, sizeof(p)) == 0) {
998
                                                                 int *r = (int *)i;
           int *r = (int *)i;
999
                                                                 r=r-1; // is this free of UB?
           *r=11; // is this free of UB?
1000
           printf("x=%d y=%d *p=%d *q=%d *r=%d\n",x,y,*p,*q,*r}r=11; // and this?
1001
                                                                 printf("x=%d y=%d *p=%d *q=%d *r=%d\n",x,y,*p,*q,*r);
         }
1002
       }
                                                               }
1003
                                                            }
```

while making the following, which tries to use the result of the cast to access both objects, UB.

```
1004
```

1005

970

- 1006 1007
- 1008

- 1014
- 1015
- 1016 1017
- 1017
- 1018
- 1020

The this, the * r=11 will resolve the provenance of the value in one way, making the r-1 UB.         In this, the * r=11 will resolve the provenance of the value in one way, making the r-1 UB.         In this, the * r=11 will resolve the provenance of the value in one way, making the r-1 UB.         In this, the * r=11 will resolve the provenance of the value in one way, making the r-1 UB.         In this, the * r=11 will resolve the provenance of the value in one way, making the r-1 UB.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance of the value in one way.         In this, the * r=11 will resolve the provenance the provenance of the value in one way. <th>1021 1022 1023 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036</th> <th><pre>// pointer_from_int_disambiguation_3.c #include <stdio.h> #include <stdin.h> #include <stdint.h> #include <inttypes.h> int y=2, x=1; int main() {     int *p = &amp;x+1;     int *q = &amp;y     uintptr_t i = (uintptr_t)p;     uintptr_t j = (uintptr_t)q;     if (memcmp(&amp;p, &amp;q, sizeof(p)) == 0) {         int *r = (int *)i;             *r=11;             r=r-1; // is this free of UB?             *r=12; // and this?             printf("x=%d y=%d *r=%d\n",x,y,*p,*q,*r);         } }</inttypes.h></stdint.h></stdin.h></stdio.h></pre></th> <th></th>	1021 1022 1023 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036	<pre>// pointer_from_int_disambiguation_3.c #include <stdio.h> #include <stdin.h> #include <stdint.h> #include <inttypes.h> int y=2, x=1; int main() {     int *p = &amp;x+1;     int *q = &amp;y     uintptr_t i = (uintptr_t)p;     uintptr_t j = (uintptr_t)q;     if (memcmp(&amp;p, &amp;q, sizeof(p)) == 0) {         int *r = (int *)i;             *r=11;             r=r-1; // is this free of UB?             *r=12; // and this?             printf("x=%d y=%d *r=%d\n",x,y,*p,*q,*r);         } }</inttypes.h></stdint.h></stdin.h></stdio.h></pre>	
In this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UE.  If this, the *r=11 will resolve the provenance of the value in one way.  If this, the *r=11 will resolve the provenance of the value in one way.  If this, the *r=11 will resolve the provenance of the value in one way.  If this, the *r=11 will resolve the provenance of the value in one way.  If this, the *r=11 will resolve the provenance of the value in one way.  If this, the *r=11 will resolve the provenance of the value in one way.  If this, the *r=11 will resolve the provenance of the value in one way.  If this, the *r=11 will resolve the provenance of the value in one way.  If this, the *r=11 will resolve the provenance of the value in one way.  If this,	1037 1038	}	
1949         1941         1943         1944         1945         1945         1947         1948         1949 <t< td=""><td>1039</td><td>In this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UB.</td><td></td></t<>	1039	In this, the *r=11 will resolve the provenance of the value in one way, making the r-1 UB.	
141         142         143         144         145         146         147         148         149         149         149         141         142         143         144         145         146         147         148         149         149         141         152         153         154         155         156         156         157         158         159         159         151         152         153         154         155         156         157         158         159         151         152         153         154         155         156         157         158         159         151         152         153         1	1040		
144         145         146         147         148         149         149         149         144         145         146         147         148         149         149         140         151         152         153         154         155         156         157         158         159         153         154         155         156         157         158         159         159         153         154         155         156         156         157         158         159         159         150         151         152         153         154         155         156         157         158         159         151         1	1041		
144         1044         1045         1046         1047         1048         1049         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1050         1051         1052         1053         1054         1055         1056         1057         1058         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1059         1051         1052         1053         1054         1055         1056 <tr< td=""><td>1042</td><td></td><td></td></tr<>	1042		
1944         1945         1947         1948         1949         1950         1951         1952         1953         1954         1953         1954         1953         1954         1955         1956         1957         1958         1959         1959         1950         1951         1952         1953         1954         1955         1956         1957         1958         1959         1959         1951         1952         1953         1954         1954         1955         1954         1955         1954         1955         1954         1955         1954         1955         1954         1955         1954         1955         1954         1955         1955 <t< td=""><td>1043</td><td></td><td></td></t<>	1043		
1045         1047         1048         1049         1051         1052         1053         1054         1055         1056         1057         1058         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1051         1052         1053         1054         1055         1056         1051         1052         1053         1054         1055         1056         1057         1058         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1059         1050         1051         1052         1053         1054 <t< td=""><td>1044</td><td></td><td></td></t<>	1044		
1046         1047         1048         1049         1050         1051         1052         1053         1054         1055         1056         1057         1058         1059         1061         1062         1053         1054         1055         1056         1057         1058         1059         1054         1055         1056         1057         1058         1059         1059         1059         1059         1059         1059         1059         1050         1051         1052         1053         1054         1055         1056         1057         1058         1059         1059         1051         1052         1053         1054         1055         1057 <t< td=""><td>1045</td><td></td><td></td></t<>	1045		
Note         1048         1049         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1061         1062         1063         1064         1065         1066         1067         1068         1069         1071         108         1072         1073         1074         1075         1076         1077         1078         1079         1074         1075         1076         1077         1078         1079         1079         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073 <tr< td=""><td>1046</td><td></td><td></td></tr<>	1046		
144         1051         1052         1053         1054         1055         1056         1057         1058         1059         1060         1061         1062         1063         1064         1065         1066         1067         1068         1079         1070         1071         1072         1073         1074         1075         1075         1076         1077         1078         1079         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075 <tr< th=""><th>1047</th><th></th><th></th></tr<>	1047		
105         106         106         106         106         106         106         107         108         107         108         109         101         102         103         1	1040		
1051         1052         1054         1055         1056         1057         1058         1059         1059         1050         1051         1052         1053         1054         1055         1056         1057         1058         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1051         1052         1053         1054         1054         1055         1056         1057         1058         1059         1059         1059         1059         1059         1059 <t< th=""><th>1019</th><th></th><th></th></t<>	1019		
1052         1053         1054         1055         1057         1058         1059         1060         1051         1052         1053         1054         1055         1056         1051         1052         1053         1054         1055         1056         1057         1058         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1059         1050         1051         1052         1053         1054         1054         1055         1056         1057         1058         1059         1059         1050         1051         1052         1053         1054         1054         1054 <t< td=""><td>1051</td><td></td><td></td></t<>	1051		
105         105         105         105         105         106         107         108         109         101         102         103         104         105         105         106         107         108         109         101         102         103         104         105         105         106         107         108         109         101         102         103         104         105         106         107         108         109         101         102         103         104         105         106         107         108         109         101         102         103         104         105         106         1	1052		
1051         1052         1053         1054         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1059         1051         1052         1053         1054         1055         1056         1057         1058         1059         1059         1051         1052         1053         1054         1054         1055         1056         1057         1058         1059         1059         1050         1051         1052         1053         1054         1054 <t< td=""><td>1053</td><td></td><td></td></t<>	1053		
105         105         105         105         106         1061         1062         1063         1064         1065         1066         1067         1068         1069         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1079         1074         1075         1076         1077         1078         1079         1074         1075         1075         1076         1077         1078         1079         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1079         1070         1070	1054		
1056         1057         1058         1059         1060         1061         1062         1063         1064         1065         1066         1067         1068         1069         1069         1070         1081         1092         1093         1094         1095         1096         1097         1098         1099         1091         1092         1093         1094         1095         1096         1097         1098         1099         1091         1092         1093         1094         1095         1096         1097         1098         1099         1090         1091         1092         1093         1094         1095         1096         1097         1098 <t< td=""><td>1055</td><td></td><td></td></t<>	1055		
1057         1058         1059         1061         1062         1063         1064         1065         1066         1067         1068         1079         1081         1092         1093         1094         1095         1096         1097         1098         1099         1091         1092         1093         1094         1095         1097         1098         1099         1091         1092         1093         1094         1095         1096         1097         1098         1099         1091         1092         1093         1094         1095         1096         1097         1098         1099         1091         1092         1093         1094         1095 <t< td=""><td>1056</td><td></td><td></td></t<>	1056		
1058         1059         1060         1061         1062         1063         1064         1065         1066         1067         1068         1070         1071         1072         1073         1074         1075         1076         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079 <t< td=""><td>1057</td><td></td><td></td></t<>	1057		
1059         1060         1061         1062         1064         1065         1066         1067         1068         1070         1071         1072         1073         1074         1075         1076         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071 <t< td=""><td>1058</td><td></td><td></td></t<>	1058		
1060         1061         1062         1063         1064         1065         1066         1067         1068         1069         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1071         1072         1073         1074         1075 <t< td=""><td>1059</td><td></td><td></td></t<>	1059		
1061         1062         1063         1064         1065         1066         1067         1068         1069         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071 <t< td=""><td>1060</td><td></td><td></td></t<>	1060		
1062         1063         1064         1065         1066         1067         1068         1069         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1071         1072         1073         1074         1075         1076         1077         1078         1079         1071         1072         1073         1074         1075         1076         1077         1078         1079         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076 <t< td=""><td>1061</td><td></td><td></td></t<>	1061		
1000         1006         1006         1007         1008         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077 <t< td=""><td>1062</td><td></td><td></td></t<>	1062		
1065 1066 1067 1068 1069 1070 1071 1072 1073 1074 1075 1076 1076 1077 1078	1064		
106         106         106         107         107         107         107         107         107         107         107         107         107         107         107         108	1065		
1067         1068         1069         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1080	1066		
1068         1069         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1079         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1071         1072         1073         1074         1075         1076         1077         1078         1079         1070         1071         1072         1073         1074         1075         1076         1077         1078         1079         1079 <t< td=""><td>1067</td><td></td><td></td></t<>	1067		
1069	1068		
1070 1071 1072 1073 1074 1075 1076 1077 1078 1079 1080	1069		
1071 1072 1073 1074 1075 1076 1077 1078 1079 1080	1070		
1072 1073 1074 1075 1076 1077 1078 1079 1080	1071		
1073 1074 1075 1076 1077 1078 1079 1080	1072		
10/4 1075 1076 1077 1078 1079 1080	1073		
1075 1076 1077 1078 1079 1080	1074		
1077 1078 1079 Draft of 1080	10/5		
1078 1079 1080	1070		
1079 1080 Draft of	1078		
1080 Draft of	1079		-
	1080		Draft of

1140

## 1081 6 TESTING THE EXAMPLE BEHAVIOUR IN CERBERUS

We have implemented executable versions of the PNVI-plain, PNVI-ae, and PNVI-ae-udi models in Cerberus [Memarian et al. 2019, 2016], closely following the detailed semantics of the accompanying note. This makes
it possible to interactively or exhaustively explore the behaviour of the examples, confirming that they are allowed
or not as intended.

01 1101 0		intended behaviour		observed behaviour				
test family	test	PNVI-plain PNVI-ae PNVI-ae-udi		Cerberus (decreasing allocator) PNVI-plain PNVI-ae PNVI-ae-udi				
	provenance_basic_global_xy.c				not triggered			
1	provenance_basic_global_yx.c		UB		UB (line 9)			
	provenance_basic_auto_xy.c	66			not triggered			
	provenance_basic_auto_yx.c			UB (line 9)				
2	cheri_03_ii.c		UB		UB (except with permissive_pointer_arith switch)			
	pointer_offset_from_ptr_subtraction_global_xy.c				UB (pointer subtraction)			
3	pointer_offset_from_ptr_subtraction_global_yx.c		UB (pointer subtraction)			Ör		
	pointer_offset_from_ptr_subtraction_auto_vy.c				UB (out-of-bou	nd store with <i>permissive_pointer_arith</i> switch		
	provenance equality global xy c					not triggered		
	provenance equality global vx.c		d			except with strict pointer equality switch)		
	provenance_equality_auto_xy.c					not triggered		
4	provenance_equality_auto_yx.c		defined, nondet		defined (ND except with strict pointer equality switch)			
	provenance_equality_global_fn_xy.c					not triggered		
	provenance_equality_global_fn_yx.c				defined (NE	D except with strict pointer equality switch)		
5	provenance_roundtrip_via_intptr_t.c		defined		defined			
	provenance_basic_using_uintptr_t_global_xy.c				not triggered			
6	provenance_basic_using_uintptr_t_global_yx.c		defined		defined			
	provenance_basic_using_uintptr_t_auto_xy.c				not triggered			
	provenance_basic_using_uintptr_t_auto_yx.c				defined			
	pointer_offset_from_int_subtraction_global_xy.c					defined		
7	pointer_offset_from_int_subtraction_global_yx.c		defined			defined		
	pointer_offset_from_int_subtraction_auto_vv.c					defined		
	pointer_offset_ron_int_subtraction_auto_yx.c	+			defined			
8	pointer_offset_xor_global.c		defined			defined		
9	provenance tag bits via uintptr t 1.c	1	defined			defined		
10	pointer_arith_algebraic_properties 2 global.c	1	defined			defined		
11	pointer_arith_algebraic_properties_3_global.c		defined			defined		
12	pointer_copy_memcpy.c		defined			defined		
13	pointer_copy_user_dataflow_direct_bytewise.c		defined			defined		
13	provenance_tag_bits_via_repr_byte_1.c		defined		defined			
15	pointer_copy_user_ctrlflow_bytewise.c		defined		defined			
16	pointer_copy_user_ctriflow_bitwise.c		defined		defined			
	provenance_equality_uintptr_t_global_xy.c				not triggered			
17	provenance_equality_unitptr_t_global_yx.c		defined			net triggered		
	provenance_equality_untptr_t_auto_xy.c				defined (true)			
	provenance union punning 2 global xv.c	defined	UB (line 16, deref)	UB (line 16, store)	a) not triggered			
	provenance union punning 2 global vx.c	defined	UB (line 16, deref)	UB (line 16, store)	defined	UB (line 16, deref) UB (line 16, store)		
18	provenance_union_punning_2_auto_xy.c	defined	UB (line 16, deref)	UB (line 16, store)		not triggered		
	provenance_union_punning_2_auto_yx.c	defined	UB (line 16, deref)	UB (line 16, store)	defined	UB (line 16, deref) UB (line 16, store)		
19	provenance_union_punning_3_global.c		defined			defined		
	provenance_via_io_percentp_global.c							
20	provenance_via_io_bytewise_global.c		filesystem	and scanf() are no	t currently sup	ported by Cerberus		
	provenance_via_io_uintptr_t_global.c	-						
	pointer_trom_integer_1pg.c	dofined (i = 7)	UB (line 7)	ino 9)	UB IN ONE EXEC (IINE 7)			
	pointer_nom_integer_ng.c		UB (Ine 6)	iiie 0)	denned (J = 7)			
	pointer from integer 1i.c	defined $(i = 7)$		ine 7)	defined $(i = 7)$			
21	pointer from integer 1ie.c		defined (i = 7)			defined (j = 7)		
	pointer_from_integer_2.c	defined (i = 7)	UB (I	ine 7)	defined (j = 7)	UB (line 7)		
	pointer from integer 2g.c		defined (j = 7)			defined (j = 7)		
	provenance_lost_escape_1.c		defined		defined			
22	provenance_roundtrip_via_intptr_t_onepast.c	UB (line 10) defined		l	JB (line 10) defined			
	pointer_from_int_disambiguation_1.c		defined ( $v = 11$ )			defined (y = 11)		
	pointer_from_int_disambiguation_1_xy.c	L	defined (y = 11)			not triggered		
23	pointer_from_int_disambiguation_2.c	LIB (line 14)		defined	UB (line 14) defined (x = 11			
	pointer_from_int_disambiguation_2_xy.c	UB (line 15)         UB (line 15)			not triggered			
	pointer_trom_int_disambiguation_3.c			UB (line 15)	UB (line 15)			
	pointer_trom_int_disambiguation_3_xy.c			not triggerea				
			areen - Cerberus ba	abaviour matches int	ent			
	(bold = tests mentioned in the document)	green - Gerberus behaviour matches intent (witch permissive pointer arith switch)						
	( <b>bold</b> = tests mentioned in the document)		hlue = Cerherus het	aviour matches inte	nt (witch nermis	Sive nointer arith Switch)		
	( <b>bold</b> = tests mentioned in the document)		blue = Cerberus ber arev = Cerberus' all	naviour matches inte ocator doesn't trigge	nt (witch <i>permis</i> the interesting	sive_pointer_arith switch)		
	(bold = tests mentioned in the document)		blue = Cerberus ber grey = Cerberus' allo	ocator doesn't trigge	nt (witch permis r the interesting	sive_pointer_anth switch) behaviour		
	( <b>bold</b> = tests mentioned in the document)		blue = Cerberus ber grey = Cerberus' allo	naviour matches inte ocator doesn't trigge	nt (witch permis r the interesting	sive_pointer_antn switch) behaviour		
	( <b>bold</b> = tests mentioned in the document)		blue = Cerberus ber grey = Cerberus' allo	naviour matches inte ocator doesn't trigge	nt (witch permis r the interesting	swe_pointer_anth switch) behaviour		
	(bold = tests mentioned in the document)		blue = Cerberus ber grey = Cerberus' allo	naviour matches inte ocator doesn't trigge	nt (witch permis r the interesting	swe_ponter_anth_switch) behaviour		
	( <b>bold = tests mentioned in the document</b> )		blue = Cerberus ber grey = Cerberus' alle	naviour matches inte ocator doesn't trigge	nt (witch permis r the interesting	swe_ponter_anth switch) behaviour		
	( <b>bold</b> = tests mentioned in the document)		blue = Cerberus bef grey = Cerberus' allo	naviour matches inte ocator doesn't trigge	nt (witch permis r the interesting	sive_pointer_anth_switch) behaviour		
	( <b>bold</b> = tests mentioned in the document)		blue = Cerberus bef grey = Cerberus' alle	naviour matches inte ocator doesn't trigge	nt (witch <i>permis</i>	sive_pointer_anth_switch) behaviour		

#### TESTING THE EXAMPLE BEHAVIOUR IN MAINSTREAM C IMPLEMENTATIONS

We have also run the examples in various existing C implementations, including GCC and Clang at various optimisation levels. 

Our test cases are typically written to illustrate a particular semantic question as concisely as possible. Some are "natural" examples, of desirable C code that one might find in the wild, but many are intentionally pathological or are corner cases, to explore just where the defined/undefined-behaviour boundary is; we are not suggesting that all these should be supported. 

Making the tests concise to illustrate semantic questions also means that most are not written to trigger inter-esting compiler behaviour, which might only occur in a larger context that permits some analysis or optimisation pass to take effect. Moreover, following the spirit of C, conventional implementations cannot and do not report all instances of undefined behaviour. Hence, only in some cases is there anything to be learned from the experimental compiler behaviour. For any executable semantics or analysis tool, on the other hand, all the tests should have instructive outcomes.

Some tests rely on address coincidences for the interesting execution; for these we sometimes include multiple variants, tuned to the allocation behaviour in the implementations we consider. Where this has not been done, some of the experimental data is not meaningful. 

The detailed data is available at https://www.cl.cam.ac.uk/~pes20/cerberus/supplementary-material-pnvi-star/ generated html pnvi star/, and summarised in the table below. 

1159					Compilers				
157		Observed behaviour (compilers), sound w.r.t PNVI-*? (relving on UB or ND?)							
160			acc-8.3		clang-7.0.1			icc-19	
100	toet family	test	PNVL-plain PNVL-ae PNV	/l-ae-udi PNVI-pla	n PNVI-20	PNVI-ae-udi	PNVI-plain	PNVI-20	PNVI-20-udi
1.71	test failing	provenance basic global xy c	v(n)	ri-ac-aai i iiivi-pia	v (n)	11111-00-001	i ivi-piain	y (y for O2+)	1100-001
1161		provenance basic global_xy.c	y (1)		pot triggorod			pot triggorod	
	1		y (y lol 02+)		not triggered			not triggered	
1162		provenance_basic_auto_xy.c	y (1)		y (11)			y (y 101 02+)	
		provenance_basic_auto_yx.c	y (h)		y (n)			y (y for 02+)	
1163	2	cheri_03_ii.c	y (n)		y (n)			y (n)	
		pointer_offset_from_ptr_subtraction_global_xy.c						y (n)	
1164	3	pointer_offset_from_ptr_subtraction_global_yx.c	y (n)		v (n)			y (n)	
1101	5	pointer_offset_from_ptr_subtraction_auto_xy.c	y (ii)		y (1)		y (y for O2+)		
11/5		pointer_offset_from_ptr_subtraction_auto_yx.c						y (y for O2+)	
1105		provenance equality global xy.c	y (n)						
		provenance equality global vx.c	y (y for Q2+)						
1166		provenance equality auto xy c	y (y for Q2+)						
	4	provenance equality auto vy c	y (n)		y (n)			y (n)	
1167		provonanco_oquality_dato_yx.c	y (n)						
		provenance_equality_global_in_xy.c	y (II)						
1168		provenance_equality_global_in_yx.c	y (y lùi O2+)						
1100	5	provenance_roundtrip_via_intptr_t.c	y (n)		y (n)			y (n)	
11/0		provenance_basic_using_uintptr_t_global_xy.c	y (n)		y (n)			n (y)	
1169	6	provenance_basic_using_uintptr_t_global_yx.c	n (y)		not triggered			not triggered	
	0	provenance_basic_using_uintptr_t_auto_xy.c	y (n)		not triggered			n (y)	
1170		provenance_basic_using_uintptr_t_auto_yx.c	y (n)		y (n)			n (y)	
		pointer_offset_from_int_subtraction_global_xy.c							
1171	-	pointer offset from int subtraction global yx.c	( )						
	/	pointer offset from int subtraction auto xy.c.	y (n)		y (n)			y (n)	
1172		pointer offset from int subtraction auto vx c							
11/2		pointer offset xor global c							
1179	8	pointer_effect_ver_globalle	y (n)		y (n)			y (n)	
11/5		provenence tog bite via viatetr t 1 e			11 (m)				
	9	provenance_tag_bits_via_unitpit_t_t.c	y (II)		y (ii)			<u>y (ii)</u>	
1174	10	pointer_aritn_aigebraic_properties_2_global.c	y (n)		y (n)			<u>y (n)</u>	
	11	pointer_arith_algebraic_properties_3_global.c	y (n)		y (n)			y (n)	
1175	12	pointer_copy_memcpy.c	y (n)		y (n)			y (n)	
	13	pointer_copy_user_dataflow_direct_bytewise.c	y (n)		y (n)			y (n)	
1176	13	provenance_tag_bits_via_repr_byte_1.c	y (n)		y (n)			y (n)	
11/0	15	pointer_copy_user_ctrlflow_bytewise.c	y (n)		y (n)			y (n)	
1177	16	pointer_copy_user_ctrlflow_bitwise.c	y (n)		y (n)			y (n)	
11//		provenance_equality_uintptr_t_global_xy.c							
4470		provenance equality uintptr t global yx.c	( )						
1178	17	provenance equality uintptr t auto xv.c	y (n) y (n)		y (n)				
		provenance equality uintptr t auto vx.c							
1179		provenance union pupping 2 global xy c	v (n)		v (n)		v (v for O2+)	n	(v)
		provenance union pupping 2 global vx c	v (v for O2+) n (v)		not triggered		, , ,	not triggered	07
1180	18	provonanco_union_punning_2_global_yx.c	y () (01 02.1)		not nggorou		w (w for O2+)	not inggered	(14)
		provenance_union_punning_2_auto_xy.c	y (1)		y (n)		y (y for O2+)		(y)
1181	10	provenance_union_punning_2_auto_yx.c	ý (II)		()		y (y 101 02+)		(y)
1101	19	provenance_union_punning_3_global.c	y (n)		y (n)			y (n)	
1100		provenance_via_io_percentp_global.c							
1182	20	provenance_via_io_bytewise_global.c	NO OPT		NO OPT			NO OPT	
		provenance_via_io_uintptr_t_global.c							
1183		pointer_from_integer_1pg.c	y (y for O0+)		y (y for O2+)			y (y for O2+)	
		pointer_from_integer_1ig.c	n (y) y (y for O2+	) n (y)	y (y f	or O2+)		n (y for O2+)	
1184		pointer_from_integer_1p.c							
		pointer from integer 1i.c	cent test with shares						
1185	21	pointer from integer 1ie.c	can't test with charon						
1100		pointer from integer 2.c							
1102		pointer from integer 2g c	v (n)		n (v)			v (n)	
1100		provenance lost escape 1 c	y (n)		v (n)			n (v for 02+)	
		provenance_lost_escape_l.c	y (n)		y (1)			11 (9 101 02.)	
1187	22	provenance_roundinp_vid_intptr_t_onepast.c	y (II)		y (II)			y (II)	
		pointer_nonl_int_disambiguation_i.c	n (y)		not utggered			not triggered	
1188		pointer_iron1_int_disampiguation_1_xy.c	not triggered		y (n)			n (y for 02+)	
	23	pointer_trom_int_disambiguation_2.c	y (n)		not triggered			not triggered	
1189		pointer_trom_int_disambiguation_2_xy.c	not triggered		y (n)			y (n)	
		pointer_trom_int_disambiguation_3.c	y (n)		not triggered			not triggered	
1100		pointer_from_int_disambiguation_3_xy.c	not triggered		y (n)			y (y for O2+)	

Page 2

(bold = tests mentioned in the document)

### 1201 REFERENCES

- 1202Mark Batty, Kayvan Memarian, Kyndylan Nienhuis, Jean Pichon-Pharabod, and Peter Sewell. 2015. The Problem of Programming Language1203Concurrency Semantics. In Programming Languages and Systems 24th European Symposium on Programming, Held as Part of the European<br/>Joint Conferences on Theory and Practice of Software, ETAPS 2015, London, UK, April 11-18, 2015. 283–307. https://doi.org/10.1007/<br/>978-3-662-46669-8\_12
- David Chisnall, Justus Matthiesen, Kayvan Memarian, Kyndylan Nienhuis, Peter Sewell, and Robert N. M. Watson. 2016. C memory object and value semantics: the space of de facto and ISO standards. http://www.cl.cam.ac.uk/~pes20/cerberus/notes30.pdf (a revison of ISO SC22
   WG14 N2013).
- Clive D. W. Feather. 2004. Indeterminate values and identical representations (DR260). Technical Report. http://www.open-std.org/jtc1/sc22/ wg14/www/docs/dr\_260.htm.
- FSF. 2018. Using the GNU Compiler Collection (GCC) / 4.7 Arrays and pointers. https://gcc.gnu.org/onlinedocs/gcc/
   Arrays-and-pointers-implementation.html. Accessed 2018-10-22.
- Krebbers and Wiedijk. 2012. N1637: Subtleties of the ANSI/ISO C standard. http://www.open-std.org/jtc1/sc22/wg14/www/docs/n1637.pdf.
   Robbert Krebbers. 2015. *The C standard formalized in Coq.* Ph.D. Dissertation. Radboud University Nijmegen.
- Juneyoung Lee, Chung-Kil Hur, Ralf Jung, Zhengyang Liu, John Regehr, and Nuno P. Lopes. 2018. Reconciling High-level Optimizations and Low-level Code with Twin Memory Allocation. In Proceedings of the 2018 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages & Applications, OOPSLA 2018, part of SPLASH 2018, Boston, MA, USA, November 4-9, 2018. ACM.
- Kayvan Memarian, Victor Gomes, and Peter Sewell. 2018. n2263: Clarifying Pointer Provenance v4. ISO WG14 http://www.open-std.org/jtc1/
   sc22/wg14/www/docs/n2263.htm.
- Kayvan Memarian, Victor B. F. Gomes, Brooks Davis, Stephen Kell, Alexander Richardson, Robert N. M. Watson, and Peter Sewell. 2019.
   Exploring C Semantics and Pointer Provenance. In *POPL 2019: Proc. 46th ACM SIGPLAN Symposium on Principles of Programming Languages.* https://doi.org/10.1145/3290380 Proc. ACM Program. Lang. 3, POPL, Article 67. Also available as ISO WG14 N2311, http://www.open-std.org/jtc1/sc22/wg14/www/docs/n2311.pdf.
- Kayvan Memarian, Justus Matthiesen, James Lingard, Kyndylan Nienhuis, David Chisnall, Robert N.M. Watson, and Peter Sewell. 2016. Into
   the depths of C: elaborating the de facto standards. In *PLDI 2016: 37th annual ACM SIGPLAN conference on Programming Language Design and Implementation (Santa Barbara).* http://www.cl.cam.ac.uk/users/pes20/cerberus/pldi16.pdf PLDI 2016 Distinguished Paper award.
- Kayvan Memarian and Peter Sewell. 2016a. N2090: Clarifying Pointer Provenance (Draft Defect Report or Proposal for C2x). ISO WG14 http://www.open-std.org/jtc1/sc22/wg14/www/docs/n2090.htm.
- Kayvan Memarian and Peter Sewell. 2016b. What is C in practice? (Cerberus survey v2): Analysis of Responses with Comments. ISO SC22
   WG14 N2015, http://www.cl.cam.ac.uk/~pes20/cerberus/analysis-2016-02-05-anon.txt.
- 1226 glibc. 2018. memcpy. https://sourceware.org/git/?p=glibc.git;a=blob;f=string/memcpy.c;hb=HEAD.
- 1227 WG14 (Ed.). 2018. Programming languages C (ISO/IEC 9899:2018 ed.).