# Refactoring for Performance with Semantic Patching: Case Study with Recipes

Michele MARTONE* and Julia LAWALL'
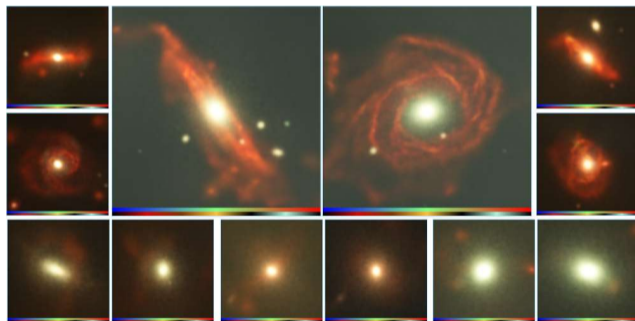
Leibniz Supercomputing Centre (Garching bei München, Germany)*
Inria, Paris, France'

C3PO'21 Workshop at ISC Frankfurt, Germany
July 2, 2021

## 66 The GADGET simulation code[1]

★ ☆ Large-scale cosmological structure formation        (galaxies and clusters)

🔧 Highly scalable        ($O(100k)$ Xeon cores on SuperMUC@LRZ)

👥 Several teams and versions        ($\geq$100 kLoC each)



*Galaxies simulated with Gadget – courtesy `http://magneticum.org`*

[1]We work with a derivative of the one described by V. Springel, "The cosmological simulation code GADGET-2", MNRAS, vol. 364, pp. 1105–1134, 2005

# **Q** Performance pilot study at LRZ

- ▶ L. Iapichino, V. Karakasis, F. Baruffa, N. Hammer[2]
- ▶ spanned over one year
- ▶ focused on 1kLoC extract
- ▶ identified changes meant for whole GADGET

**"**

# 🔨 Speedup requires data layout change ⇄

**main header:**

```
1 struct particle {
2   double  Mass, ...
3 #if defined(BLACK_HOLES)...
4   double  Hsml, ...
5   ...
6 };
```

$\Longrightarrow$

```
1 struct particle_soa_t {
2   double *Mass, ...
3 #if defined(BLACK_HOLES)...
4   double *Hsml, ...
5   ...
6 };
```

**init source file:**

```
1 // Array of Structures:
2 struct particle *P;
3 // allocate one global array:
4 P = mymalloc(...
5
6
7
8
```

$\Longrightarrow$

```
1 // Structure of Arrays
2 struct particle_soa_t P_SoA;
3 // allocate one global array
4 // for each field:
5 P_SoA.Mass = mymalloc(...
6 #if defined(BLACK_HOLES)...
7 P_SoA.Hsml = mymalloc(...
8 ...
```

**favour auto-vectorization in ∗.c:**

```
1 ...
2 // may not vectorize
3 P[i].Mass + P[i]...
```

$\Longrightarrow$

```
1 ...
2 // vectorizes better
3 P_SoA.Mass[i] + P_SoA...
```

# ✺ Speedup *suggests* major code change steps ⛶

## 1. each type, $10 \approx 100$ fields, `#ifdefs`:

```
1 struct particle {
2   double  Mass, ...
3 #if defined(BLACK_HOLES)...
4   double  Hsml, ...
5   ...
6 };
```

$\Longrightarrow$

```
1 struct particle_soa_t {
2   double *Mass, ...
3 #if defined(BLACK_HOLES)...
4   double *Hsml, ...
5   ...
6 };
```

## 2. (almost) each field an allocation:

```
1 // Array of Structures:
2 struct particle *P;
3 // allocate one global array:
4 P = mymalloc(...
5
6
7
8
```

$\Longrightarrow$

```
1 // Structure of Arrays
2 struct particle_soa_t P_SoA;
3 // allocate one global array
4 // for each field:
5 P_SoA.Mass = mymalloc(...
6 #if defined(BLACK_HOLES)...
7 P_SoA.Hsml = mymalloc(...
8 ...
```

## 3. ≫10KLoC change in `*.c`!

```
1 ...
2 // may not vectorize
3 P[i].Mass + P[i]...
```

$\Longrightarrow$

```
1 ...
2 // vectorizes better
3 P_SoA.Mass[i] + P_SoA...
```

# Tool choice criteria 🔧

1. *timeliness*: How to quickly change so many lines of code?
2. *correctness*: How to avoid introducing mistakes? Note the aggravation of having numerous build-time *code paths* implied by the `#ifdefs`.
3. *flexibility*: Can we enact only a *partial* SoA translation, possibly on demand?
4. *continuity*: Can we develop in AoS, transforming only before *build and run*?
5. *acceptance*: How to have the community *accept* the proposed solution?

# A C code matching and transformation engine ⚙

▶ A project from Inria (France)



---

# A C code matching and transformation engine ⚙

▶ A project from Inria (France)

▶ originally to

    🐧 update Linux kernel drivers[3]

    🐞 **smash bugs**
    (hence the name)[4]



---

[3]https://git.kernel.org/pub/scm/linux/kernel/git/backports/backports.git/tree/patches
[4]https:
//git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/tree/scripts/coccinelle

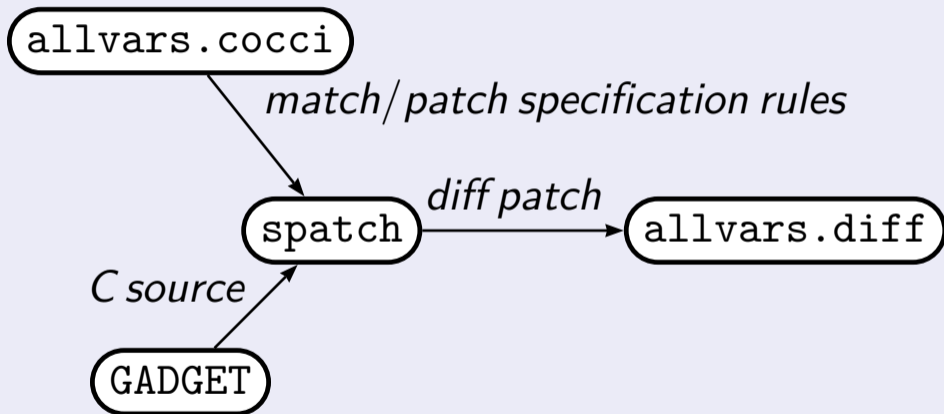# A C code matching and transformation engine ⚙

▶ A project from Inria (France)

▶ originally to

    🐧 update Linux kernel drivers[3]

    🐛 **smash bugs**
    (hence the name)[4]
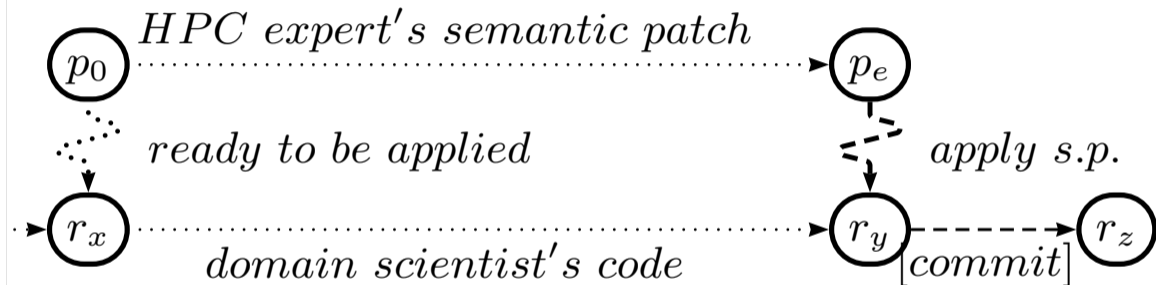
▶ seemingly **underutilized** in other contexts



---

[3]https://git.kernel.org/pub/scm/linux/kernel/git/backports/backports.git/tree/patches
[4]https:
//git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/tree/scripts/coccinelle
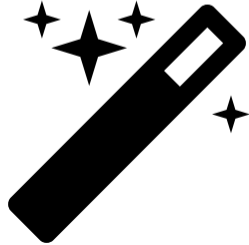
## Patch what and how ?

# Transformations *encoded* in *semantic patches* 🔧



- produces usable patched (e.g. refactored) code
- no requirement to commit patched code

# Step 1: modify data structures into SoA

# Match main data structure 🧩

```
1 @prtcl_str@
2 identifier id = {particle_data,sph_particle_data};
3 field list fs;
4 identifier I;
5 @@
6
7 struct id { fs } *I;
```

- ▶ pattern language to *match* C
- ▶ *metavariables* `id`, `I`, ... to match C language elements

# Create new identifier 🏷️

```
1 @script:python new_prtcl_str_id@
2 id << prtcl_str.id;
3 id1;
4 @@
5
6 coccinelle.id1="%s_soa_t"%(id)
```

- ▶ reuse `identifier` `id` from rule `prtcl_str`
- ▶ new string using Python

# Clone main data structure 📑

```
1 @insert_new_prtcl_str depends on prtcl_str@
2 identifier new_prtcl_str_id.id1;
3 field list prtcl_str.fs;
4 @@
5
6 extern int maxThreads;
7 ++struct id1 { fs };
```

- ▶ pattern language to *patch* C
- ▶ reminiscent of GNU patch
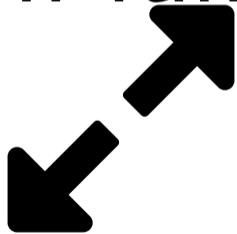
# Whitelisted types to pointers ♻

```
1  @make_ptr@
2  identifier new_prtcl_str_id.id1;
3  identifier M;
4  typedef MyDouble;
5  typedef MyFloat;
6  typedef MyLongDouble;
7  typedef MyDoublePos;
8  typedef MyBigFloat;
9  type MT = { double,float,MyDouble,MyFloat,MyLongDouble,MyDoublePos,MyBigFloat};
10 @@
11
12 struct id1 { ...
13 - MT   M;
14 + MT *M;
15    ...
16 };
```

# Remove non-pointer fields ✖

```
1 @del_non_ptr@
2 identifier new_prtcl_str_id.id1;
3 identifier J;
4 type T;
5 type P != {T*};
6 @@
7
8 struct id1 { ...
9 - P J;
10   ...
11 };
```

▶ *for each field* $J$ *of type* $P$ *which is <u>not</u> a pointer to another type...*

# Step 2: allocation functions

## Populate allocation functions ↗

```
1 @per_type_soa_alloc@
2 identifier new_prtcl_str_id.id1;
3 identifier prtcl_str_mmbrs.M;
4 type prtcl_str_mmbrs.MT;
5 symbol P;
6 identifier insert_per_type_soa_functions.soa_alloc_fid;
7 fresh identifier si = "#ifdef "##"HAVE__"##id1##"__"##M##" //";
8 identifier str_from_id.str;
9 @@
10
11 void soa_alloc_fid(...)
12 { ...
13 ++si;
14 ++      P->M = (MT*) mymalloc(str, sizeof(*(P->M)) * N);
15 ++      if(!P->M)soa_abort(/*"allocating "*/ str);
16 ++#endif
17 }
```

# Step 3: AoS ➜ SoA ⚡

# Match old structure defs[5] 回

```
1 @ostr@
2 identifier id = {particle_data,sph_particle_data};
3 identifier P;
4 @@
5
6 struct id { ... } *P;
```

---
[5]Subset of rule prtcl_str from p. 12.

# Match fields of new types[6] 📌

```
1 @nt@
2 identifier pps.id1;
3 identifier I;
4 type T;
5 @@
6
7 struct  id1 {
8  ...
9  T I;
10  ...
11 };
```

---

[6]Relies on rule pps, duplicate of `new_prtcl_str_id`, p. 13.

# Create new identifiers[7] 🏷️

```
1 @script:python pid@
2 id1 << pps.id1;
3 P << ostr.P;
4 S;
5 @@
6
7 coccinelle.S="%s_soa"%(P)
```

---

[7]Same mechanism as `new_prtcl_str_id` at p. <span style="color:red">13</span>.

# Patch many expressions (≫10KLoC diff)

```
1 @soa_access@
2 identifier ostr.P;
3 identifier pid.S;
4 identifier nt.I;
5 expression E;
6 @@
7
8 - P[E].I
9 + S.I[E]
```

∀ `identifier` P previously matched in rule `ostr`
∀ `identifier` S previously matched in rule `pid`
∀ `identifier` I previously matched in rule `nt`
∀ `expression` E matching rule `soa_access`

substitute `P[E].I` with `S.I[E]`

# Lessons learned

Easier semantic patching if ..

- ▶ Adopt *coding guidelines*
- ▶ Follow emerging practices in *research software engineering*

# Outcome ✔

🚀 GADGET can evolve further

▶ its semantic patches...

    ▶ ...can be stored 💾

    ▶ ...and applied at anytime ▶

    ▶ ...ease performance experiments ⚗

    ▶ ...serve also as a Coccinelle test[8] ♻

---

[8] https://github.com/coccinelle/coccinelle/commit/ad5a94