

# PLANKTON: Reconciling Binary Code and Debug Information (Supplement Material)

August 11, 2023

## 1. Formalization of Transformation Rules

We follow existing formalization methods [3, 1, 2] for LLVM intermediate representations. Figure 1 shows (a fragment of) the abstract syntax for the subset of the LLVM IR formalized in PLANKTON, which is used to define peephole transformation rules. Operations in the LLVM IR compute with values  $op$ , which are either identifiers  $id$  naming temporaries, or constants computed from statically-known data.

Figure 2 shows more transformation rules. `subToAdd` and `ptrToPtr` are two auxiliary transformation rules. `subToAdd` converts pointer subtraction to pointer addition, which allows PLANKTON to apply `multiAdd`. `ptrToPtr` simplifies redundant type conversion produced by other rules. `ptrIcmp` and `ptrPhi` transform values that are used in `phi`-nodes and comparisons, where all operands should have consistent types. These two rules use type hints of one incoming value to transform all the others. `loadBitcast`, `storeBitcast`, `loadPtr`, and `storePtr` are four rules used to find correct types for memory load/store operations. `promoteLoad` and `promotePtrToInt` aggregate type hints based on all instructions that use the value and apply those hints to the current instruction.

The above promotion rules assign values with correct types and propagate type hints through the control- and data-flows because type conversions are passed on to the following instructions and uses of the values. The demotion rules are just a set of “inverse” rules for the promotion. For example, `intToPtrGep` will lower the `getelementptr` instruction into a sequence of `add` instructions.

## References

- [1] Juneyoung Lee, Chung-Kil Hur, Ralf Jung, Zhengyang Liu, John Regehr, and Nuno P Lopes. Reconciling high-level optimizations and low-level code in llvm. *Proceedings of the ACM on Programming Languages*, 2(OOPSLA):1–28, 2018.
- [2] Nuno P Lopes, David Menendez, Santosh Nagarakatte, and John Regehr. Provably correct peephole optimizations with alive. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 22–32, 2015.
- [3] Jianzhou Zhao, Santosh Nagarakatte, Milo MK Martin, and Steve Zdancewic. Formalizing the llvm intermediate representation for verified program transformations. In *Proceedings of the 39th annual ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 427–440, 2012.

```

prog      ::=  $\frac{pre\ nl\ stmt}{stmt}$ 
stmt      ::=  $stmt\ nl\ stmt \mid reg = inst \mid inst$ 
cond      ::= eq | ne | ugt | uge | ult | ule | sgt | sge | slt | sle
typ       ::= isz | void |  $typ* \mid \{\overline{typ}_j^j\}$ 
binop     ::= add | sub | mul | udiv | sdiv | urem | srem
           | shl | lshr | ashr | and | or | xor
castop    ::= ptrtoint | inttoptr | bitcast
 $\phi$       ::= phi  $typ \llbracket [op_j, b_j] \rrbracket_j$ 
inst      ::=  $binop\ typ\ op_1, op_2 \mid castop\ typ_1\ op_1\ to\ typ_2\ op_2$ 
           | load  $typ* op \mid store\ typ\ op_1\ typ* op_2$ 
           | select  $op_1\ typ\ op_2\ op_3 \mid icmp\ cond\ typ\ op_1, op_2$ 
           | getelementptr  $typ* op\ isz \llbracket [op_j] \rrbracket_j$ 

```

Figure 1: Partial syntax for peephole transformation rules.

Table 1: Bug detection results for IRs produced by WLLVM and lifted from GCC compiled binaries by PLANKTON. ▼ means the static analysis runs out of memory. ▲ means lifting failure. Only the client sides of MySQL and MariaDB are invoked since both SVF and PINPOINT timeout on their server sides.

Program	KLOC	WLLVM		PLANKTON					
		PINPOINT # Report	SVF # Report	PINPOINT			SVF		
				#FP	#FN	#TP	#FP	#FN	#TP
CWE401s01	188	321	764	83	8	313	55	117	647
CWE401s02		387	0	77	32	355	0	0	0
CWE401s03		193	450	53	6	187	36	78	372
CWE415s01	103	446	361	28	120	326	15	191	170
CWE415s02		302	0	32	84	218	0	0	0
CWE416	57	884	236	60	1	883	0	0	236
CWE476	36	280	51	10	150	130	0	0	51
bash	397	3	31	12	1	2	60	7	24
darknet	28	21	666	5	6	15	101	138	528
ffmpeg	1162	196	▼	211	90	106	▼	▼	▼
git	262	68	612	47	48	20	161	79	533
libcrypto.so.3	490	114	647	21	41	73	418	50	597
libcuc.so	833	97	257	66	34	63	34	71	186
libuv.so.1.0.0	62	0	7	1	0	0	0	0	7
mariadb	3966	39	86	2	11	28	15	13	73
mysql	4572	64	111	6	26	38	20	19	92
php	1358	189	2160	51	84	105	574	206	1783
python	749	31	719	7	15	16	326	421	420
redis-server	170	871	911	214	387	487	93	697	214
ss-local	34	1	22	4	1	0	4	6	16
tmux	60	49	292	35	10	39	43	116	176
vim	389	103	536	13	33	70	131	111	425
wget	129	42	53	19	4	38	25	12	41
wrk	596	56	516	48	33	23	217	100	416
<b>Total</b>		<b>4757</b>	<b>9488</b>	<b>1105</b>	<b>1225</b>	<b>3535</b>	<b>2328</b>	<b>2432</b>	<b>7007</b>
<b>Rate</b>				<b>23.81%</b>	<b>25.74%</b>	<b>74.26%</b>	<b>24.94%</b>	<b>25.77%</b>	<b>74.23%</b>

**Promotion Rules**

$$\begin{array}{c}
\frac{v = \text{ptrtoint } \tau^* v_p \text{ to isz} \quad (isAgg(\tau) \wedge (v_k = v \vee v_j = v)) \Rightarrow \llbracket v_i = \text{add isz } v_j, v_k \rrbracket_n}{v' = \text{getelementptr } \tau^* v_p, \text{ i32 } \llbracket v_{idx} \rrbracket_m \quad v = \text{ptrtoint typeOf}(v') \quad v' \text{ to isz}} \quad (\text{multiAdd}) \\
\\
\frac{v_i = \text{ptrtoint } \tau^* v_j \text{ to isz} \quad v = \text{inttoptr isz } v_i \text{ to } \tau^*}{\text{replaceUse}(v, v_j)} \quad (\text{ptrToInt}) \quad \frac{\exists k, v_k = \text{ptrtoint } \tau^* v_p \text{ to isz}; v = \text{phi isz } \llbracket [v_i, b_i] \rrbracket_n}{\begin{array}{l} \llbracket v'_i = \text{typeConv}(v_i, \tau^*) \rrbracket_n \\ v_i = \text{phi } \tau^* \llbracket [v'_i, b_i] \rrbracket_n \\ v = \text{ptrtoint } \tau^* v_i \text{ to isz} \end{array}} \quad (\text{ptrPhi}) \\
\\
\frac{v'_p = \text{bitcast } ty_1^* v_p \text{ to } ty_2^* \quad v = \text{load } ty_2^* v'_p \quad \phi \subset \{ty_1 > ty_2\}}{v'_p = \begin{cases} \text{getelementptr } ty_1^* v_p, 0 & \text{isAgg}(ty_1) \\ v_p & \text{otherwise} \end{cases} \quad v' = \text{load typeOf}(v'_p) v'_p \quad v = \text{typeConv}(v', ty_2)} \quad (\text{loadBitcast}) \quad \frac{v'_p = \text{bitcast } ty_1^* v_p \text{ to } ty_2^* \quad \text{store } ty_2 v, ty_2^* v'_p \quad \phi \subset \{ty_1 > ty_2\}}{v'_p = \begin{cases} \text{getelementptr } ty_1^* v_p, 0 & \text{isAgg}(ty_1) \\ v_p & \text{otherwise} \end{cases} \quad v' = \text{typeConv}(v, \text{typeOfPointee}(v'_p)) \quad \text{store typeOf}(v') v', \text{typeOf}(v'_p) v'_p} \quad (\text{storeBitcast}) \\
\\
\frac{\llbracket v_1 \mid v_2 \rrbracket = \text{ptrtoint } \tau^* v_p \text{ to isz} \quad \text{icmp isz } v_1, v_2}{v'_1 = \text{typeConv}(v_1, \tau^*) \quad v'_2 = \text{typeConv}(v_2, \tau^*) \quad \text{icmp } \tau^* v'_1, v'_2} \quad (\text{ptrIcmp}) \quad \frac{v'_p = \text{load isz}^* v_p; \quad v = \text{inttoptr isz } v'_p \text{ to } \tau^*}{v'_p = \text{bitcast isz}^* v_p \text{ to } \tau^{**} \quad v_l = \text{load } \tau^{**} v'_p \quad v' = \text{ptrtoint } \tau^* v_l \text{ to isz} \quad v = \text{inttoptr isz } v' \text{ to } \tau^*} \quad (\text{loadPtr}) \quad \frac{v'_p = \text{ptrtoint } \tau^* v_p \text{ to isz} \quad \text{store isz } v'_p, \text{ isz}^* v}{v' = \text{bitcast isz}^* v \text{ to } \tau^{**}; \quad \text{store } \tau^* v_p, \tau^{**} v'} \quad (\text{storePtr}) \\
\\
\frac{v = \text{load isz}^* v_p \quad (n = 0 \wedge (v_k = v \vee v_j = v)) \Rightarrow \llbracket v_i = \text{add isz } v_j, v_k \rrbracket_n \quad \exists m, v_m = \text{inttoptr isz } v \text{ to } \tau^*}{v = \text{load isz}^* v_p \quad v' = \text{inttoptr isz } v \text{ to } \tau^* \quad v'' = \text{ptrtoint } \tau^* v' \text{ to isz} \quad (n = 0 \wedge (v_k = v'' \vee v_j = v'')) \Rightarrow \llbracket v_i = \text{add isz } v_j, v_k \rrbracket_n} \quad (\text{promoteLoad}) \quad \frac{v = \text{ptrtoint } ty_1^* v_p \text{ to isz} \quad \exists i, v_i = \text{inttoptr isz } v \text{ to } ty_2^* \quad \phi \subset \{ty_2 > ty_1\}}{v_i = \text{bitcast } ty_1^* v_p \text{ to } ty_2^* \quad v_j = \text{ptrtoint } ty_2^* v_i \text{ to isz} \quad \text{replaceUse}(v, v_j)} \quad (\text{promotePtrToInt}) \\
\\
\text{Demotion Rules} \\
\\
\frac{\exists k, v_k = \text{inttoptr isz } v_p \text{ to } \tau^*; v = \text{phi } \tau^* \llbracket [v_i, b_i] \rrbracket_n}{\begin{array}{l} \llbracket v'_i = \text{typeConv}(v_i, \text{isz}) \rrbracket_n \\ v_i = \text{phi isz } \llbracket [v'_i, b_i] \rrbracket_n \\ v = \text{inttoptr isz } v_i \text{ to } \tau^* \end{array}} \quad (\text{intPhi}) \quad \frac{\llbracket v_1 \mid v_2 \rrbracket = \text{inttoptr isz } v_b \text{ to } \tau^* \quad \text{icmp } \tau^* v_1, v_2}{v'_1 = \text{typeConv}(v_1, \text{isz}) \quad v'_2 = \text{typeConv}(v_2, \text{isz}) \quad \text{icmp isz } v'_1, v'_2} \quad (\text{intIcmp}) \\
\\
\frac{v_p = \text{inttoptr isz } v_b \text{ to } \tau^* \quad v = \text{getelementptr } \tau^* v_p, \text{ i32 } \llbracket v_{idx} \rrbracket_m \quad \tau' = \text{typeOf}(v)}{v_p = \text{inttoptr isz } v_b \text{ to } \tau^* \quad v_b = \text{ptrtoint } \tau^* v_p \text{ to isz} \quad (n = 0 \wedge (v_k = v_b \vee v_j = v_b)) \Rightarrow \llbracket v_i = \text{add isz } v_j, v_k \rrbracket_n \Rightarrow v_n \quad v = \text{inttoptr isz } v_n \text{ to } \tau'} \quad (\text{intToPtrGep}) \quad \frac{v'_p = \text{load } \tau^{**} v_p; \quad v = \text{ptrtoint } \tau^* v'_p \text{ to isz}}{v'_p = \text{bitcast } \tau^{**} v_p \text{ to } \tau^* \quad v_l = \text{load } \tau^* v'_p \quad v' = \text{inttoptr isz } v_l \text{ to } \tau^* \quad v = \text{ptrtoint } \tau^* v_l \text{ to isz}} \quad (\text{loadInt}) \quad \frac{v'_p = \text{inttoptr isz } v_p \text{ to } \tau^* \quad v = \text{ptrtoint } \tau^* v'_p \text{ to isz}}{\text{replaceUse}(v, v_p)} \quad (\text{intToPtr})
\end{array}$$

Figure 2: Selected promotion and demotion transformation rules.