

LIBLISA

Instruction Discovery and Analysis on x86-64

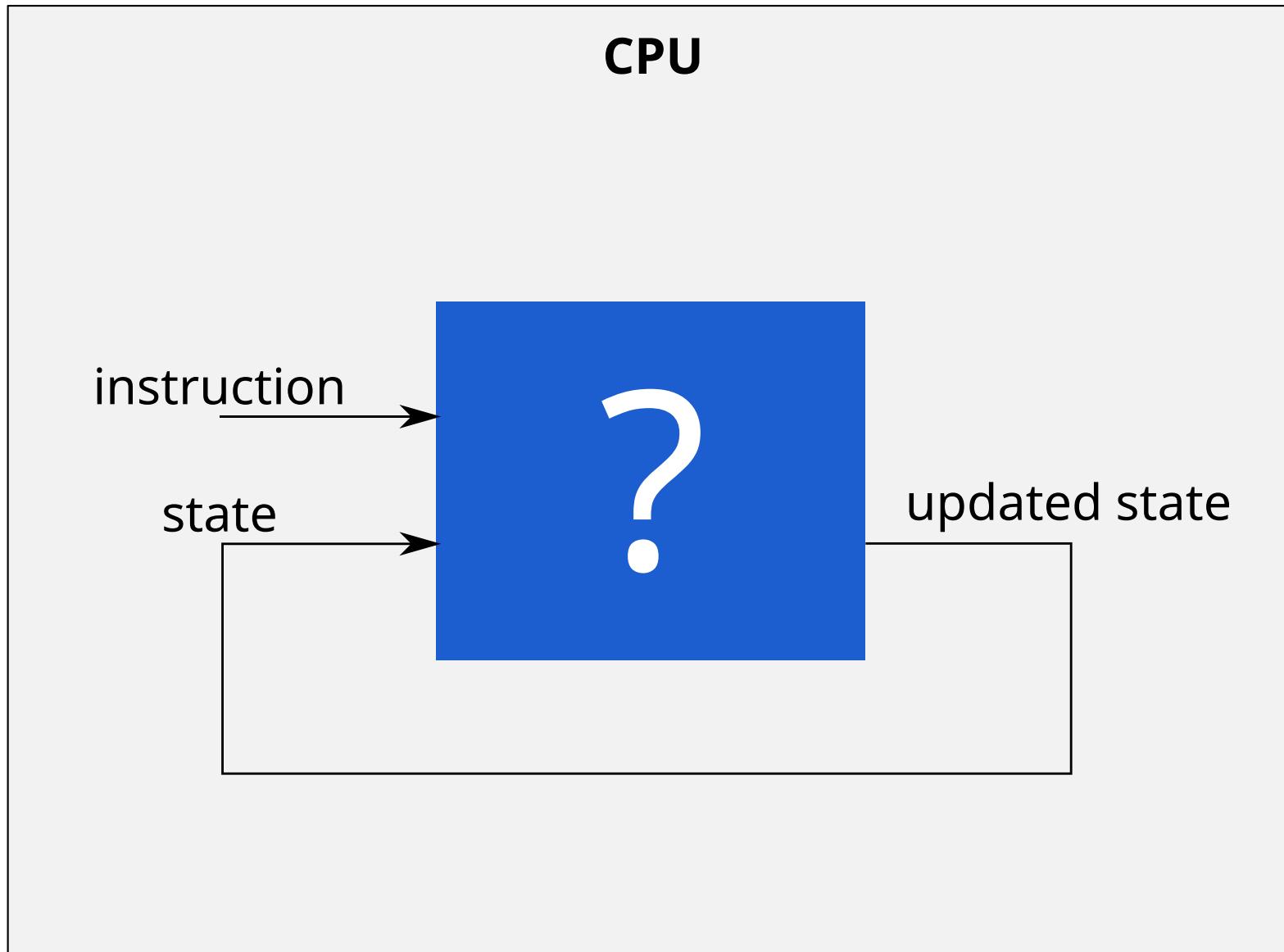


Open Universiteit



VIRGINIA
TECH

CPU Execution Model



Instruction Semantics

48 89 f8



Instruction Semantics



- $r_{RDI} := r_{RAX}$
- $r_{RIP} := r_{RIP} + 3$

x86 Instruction Semantics

A screenshot of a PDF viewer window titled "Intel® 64 and IA-32 Architectures Software Developer's Manual Combined Volumes 1, 2A, 2B, 2C, 2D, 3A, 3B, 3C, 3D, and 4 — Okular". The window shows the title page of the manual. A large white circle highlights the text "of 5082" in the top right corner. Below the title, the Intel logo is visible. The title text is in large blue font, and the volume information is in smaller blue font.

Combined Volumes: 1, 2A, 2B, 2C, 2D, 3A, 3B, 3C, 3D, and 4

of 5082

Intel® 64 and IA-32 Architectures Software Developer's Manual

Combined Volumes:
1, 2A, 2B, 2C, 2D, 3A, 3B, 3C, 3D, and 4

NOTE: This document contains all four volumes of the Intel 64 and IA-32 Architectures Software Developer's Manual: *Basic Architecture*. Order Number 253665: *Instruction Set Reference A-Z*. Order

Instruction semantics are hard

Inconsistencies Found in the Intel Manual

Errors in Intel documentation The Intel manual only contains informal pseudo-code and an English description of the semantics, so there is no automatic way of comparing the learned formulas with the manual. However, over several inconsistencies

Errors in Intel documentation

the over 3,800 pages of the manual, the majority of which is not machine-checked. Examples of these errors include:

- The upper 128 bits of the output register of `pshufb xmm0, xm` are kept unmodified on Intel hardware, but the pseudo-code in the manual incorrectly says they are cleared.
- The *Op/En* column for the `vextracti128` instruction should be MRI, but it was RMI. This causes a wrong operand ordering when assembling this instruction.
- One possible encoding for `xchgl eax, ea` is `0x90 0x00` according to the manual, and its semantics are to clear the upper 32 bits of `rax`. However, `0x90` is also the encoding of `nop`, which does nothing according to the manual (and hardware).
- The textual description of `roundpd` claims to produce a single-precision floating point value, when it produces a double-precision floating point value.

We have reported all of these items and they were confirmed by Intel as errors in the manual [8].

Inconsistencies Found in the Intel Manual Here are inconsistencies found during development and testing. According to the manual, the semantics of `vpsravd %xmm3, %xmm2, %xmm1` seems to depend on the lower 100 bits of `%xmm3`, whereas the actual hardware execution suggests that it should depend on the lower 128 bits. Similar inconsistencies are found in instructions with mnemonics `vpslld`, `vpsllvq`, `vpsravd`. Also, we found misleading typos related to instructions with opcodes `vpsravw`, `vpsravd`, `vpsravd`.

observed behaviors contradicts the x86 reference manual

We also discovered cases where observed behaviors contradicts the x86 reference manual (which is unsurprising given the size and complexity of the spec). For instance, we discovered while debugging our template *T-ARI^{main}* that the overflow OF flag should be set to 0 after executing IMUL8 with 65 and 254 as inputs according to the Intel spec, while the OF flag is actually set to 1 after the execution of this instruction with those inputs on an Intel XEON3.7 processor.

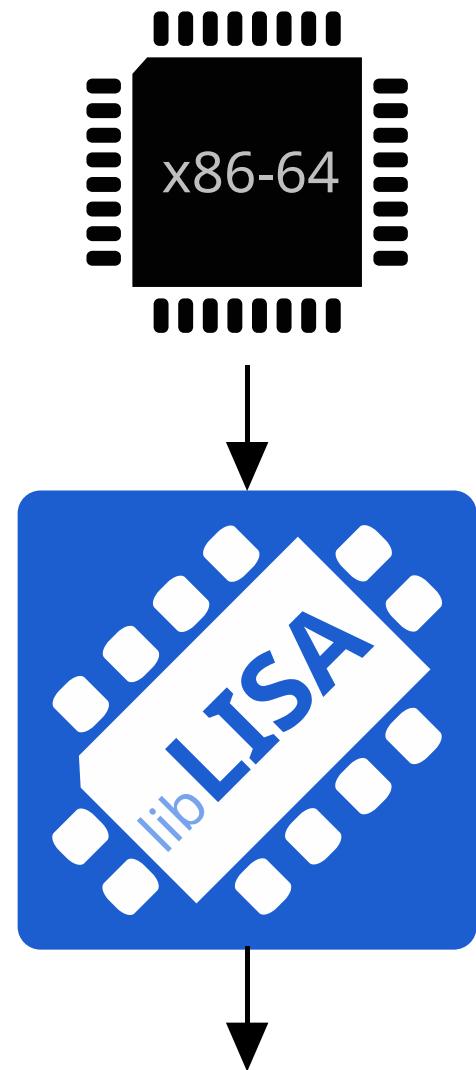
Dasgupta et al., "A complete formal semantics of x86-64 user-level instruction set architecture" (PLDI'19);

Heule et al., "Stratified Synthesis: Automatically Learning the x86-64 Instruction Set" (PLDI'16);

Godefroid and Taly, "Automated Synthesis of Symbolic Instruction Encodings from I/O Samples" (PLDI'12)

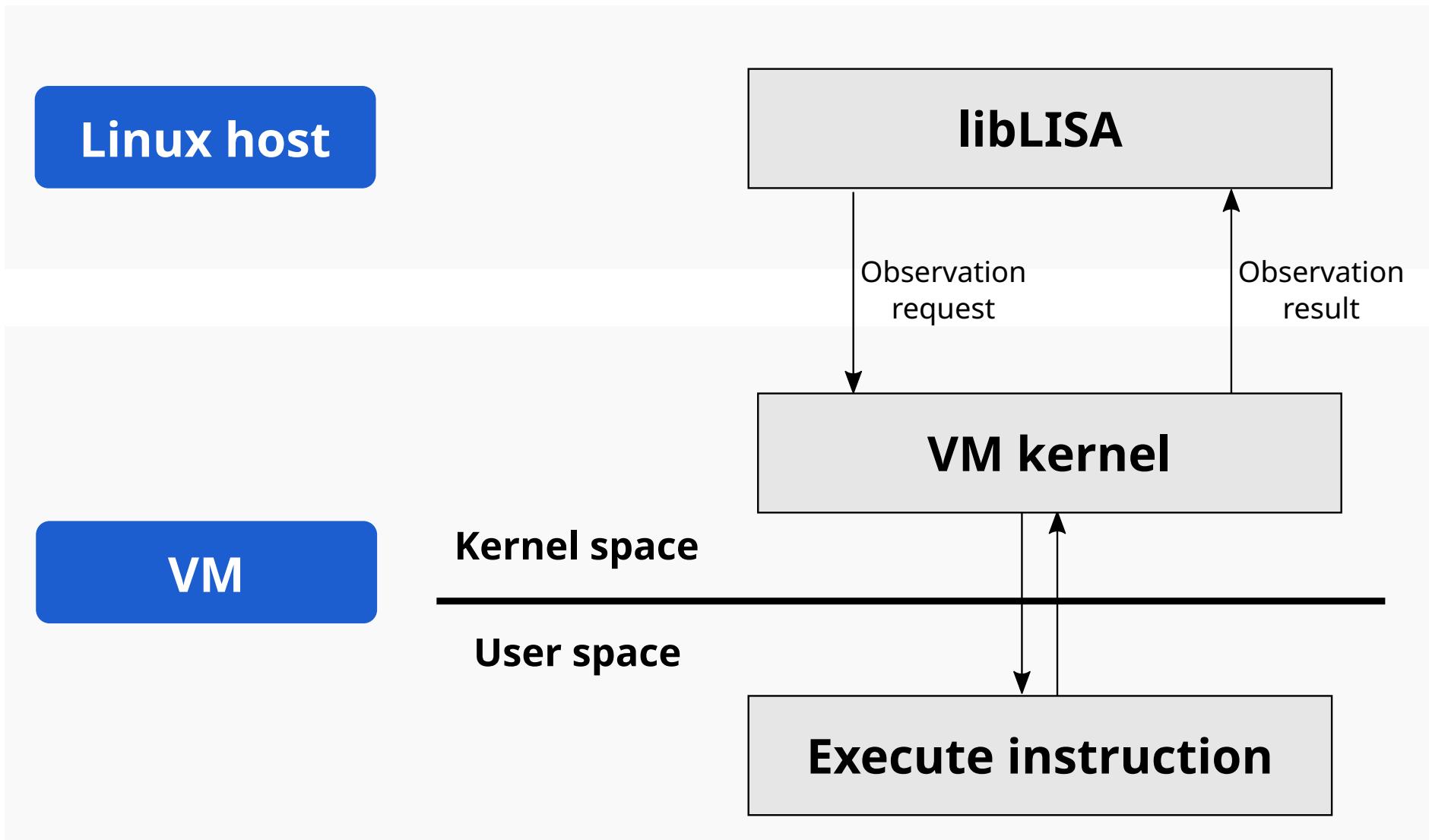
We should automate this

- **Specify only the essentials:**
 - Description of CPU state
 - CPU Observer
- **Automatically infer the rest**



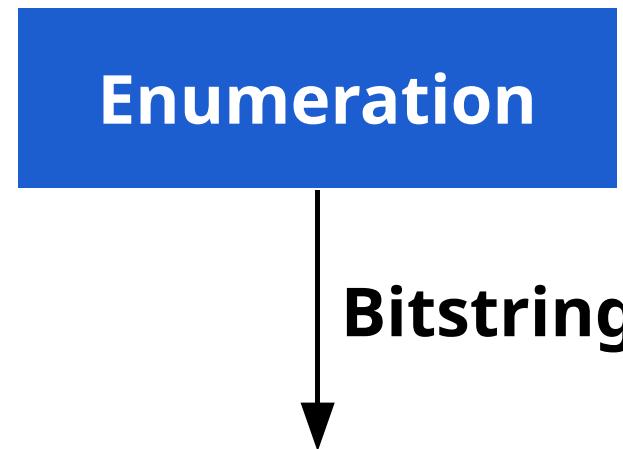
CPU semantics

Instruction Execution Observation

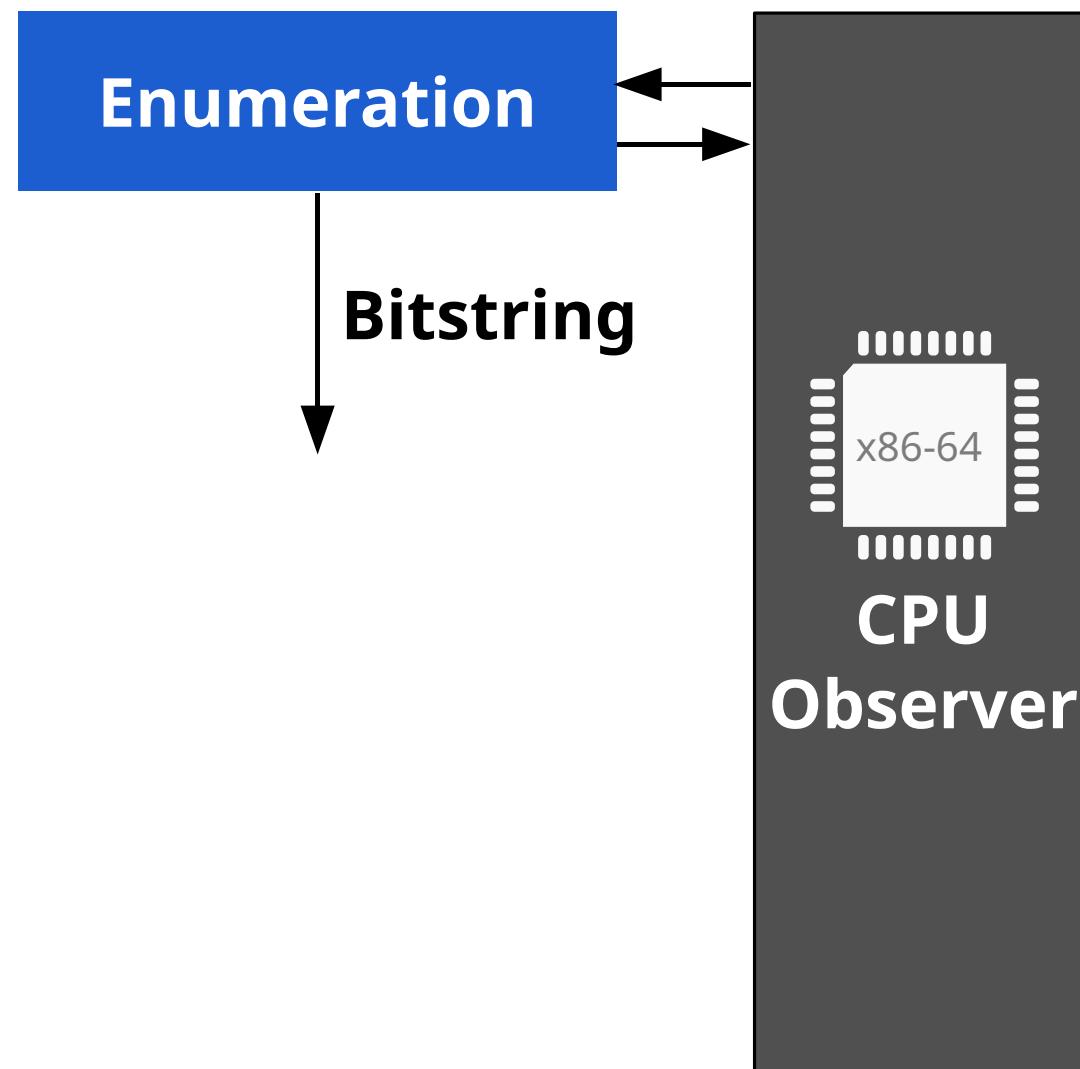


How do we automate everything?

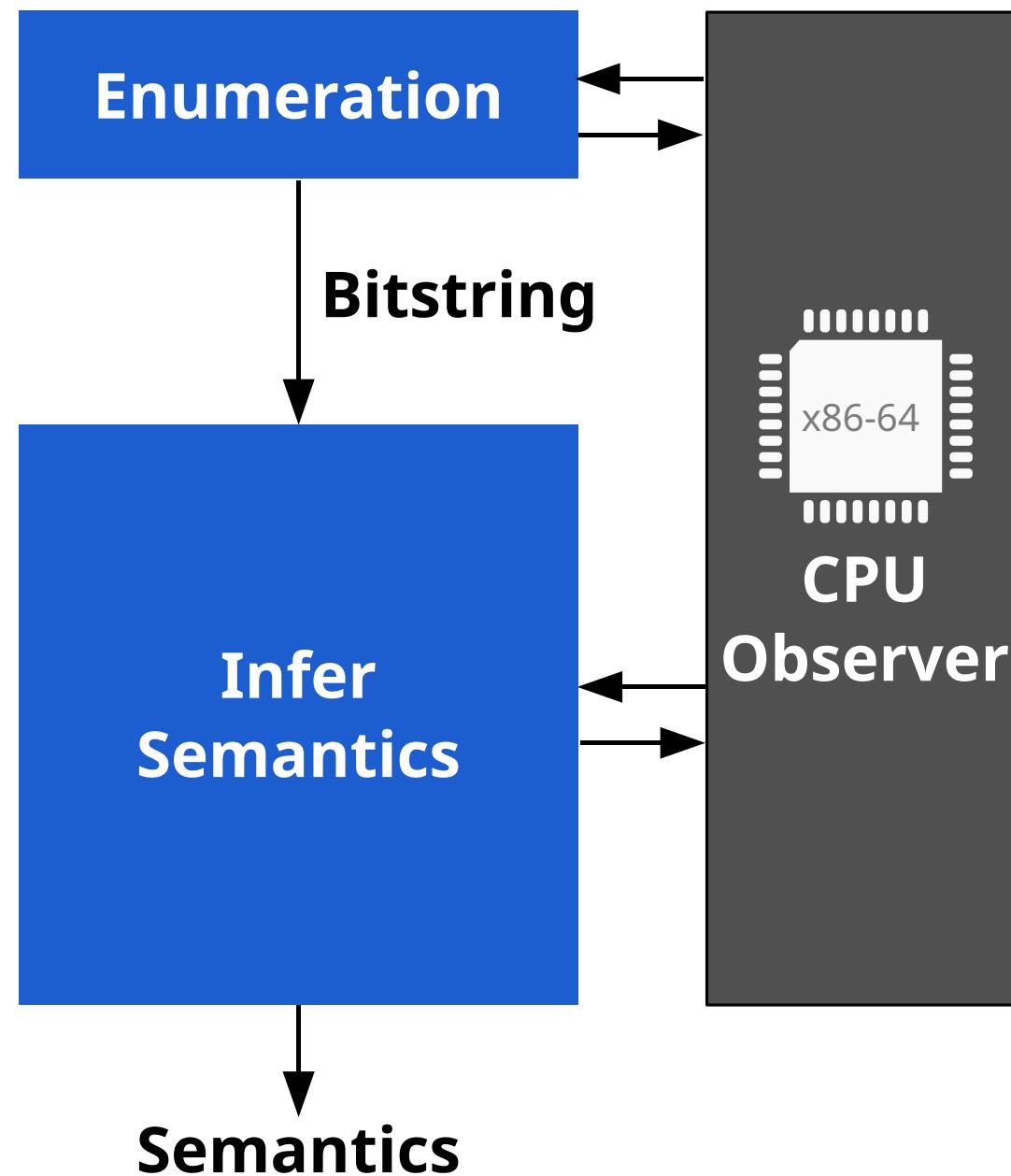
How do we automate everything?



How do we automate everything?



How do we automate everything?



Dealing with long instructions

10111000 01111000 01010110 00110100 00010010

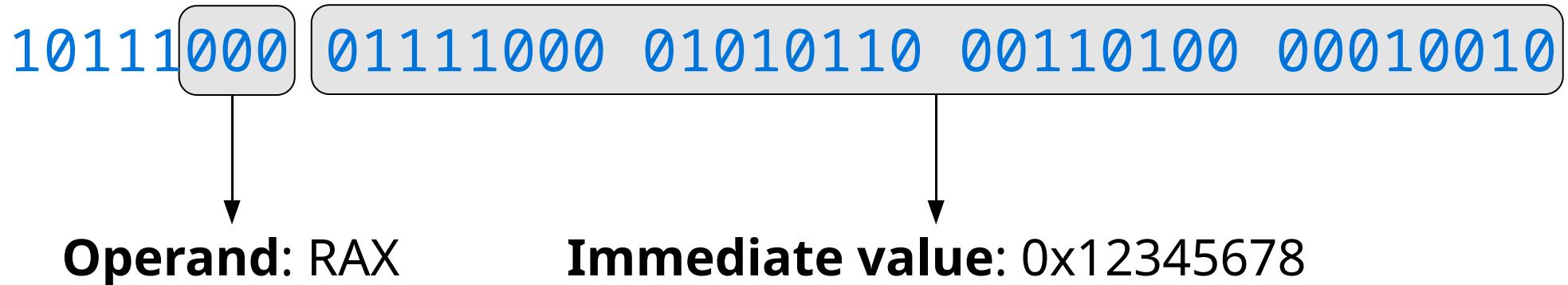
Dealing with long instructions

10111000 01111000 01010110 00110100 00010010

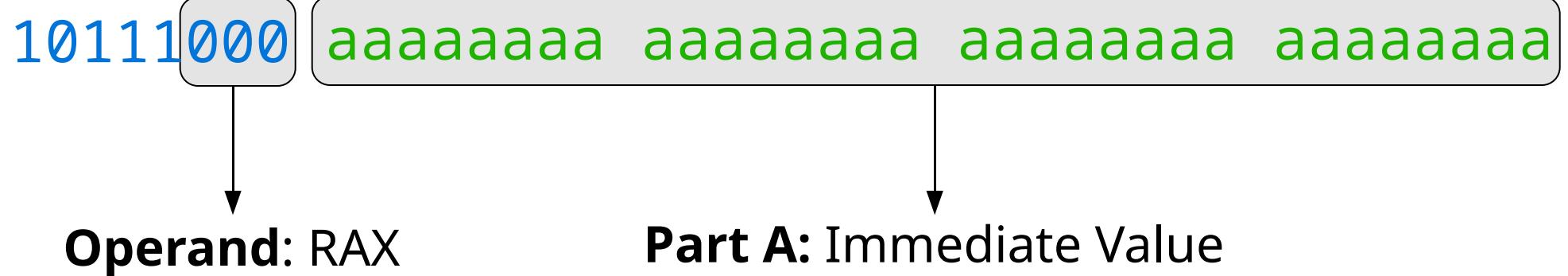


Immediate value: 0x12345678

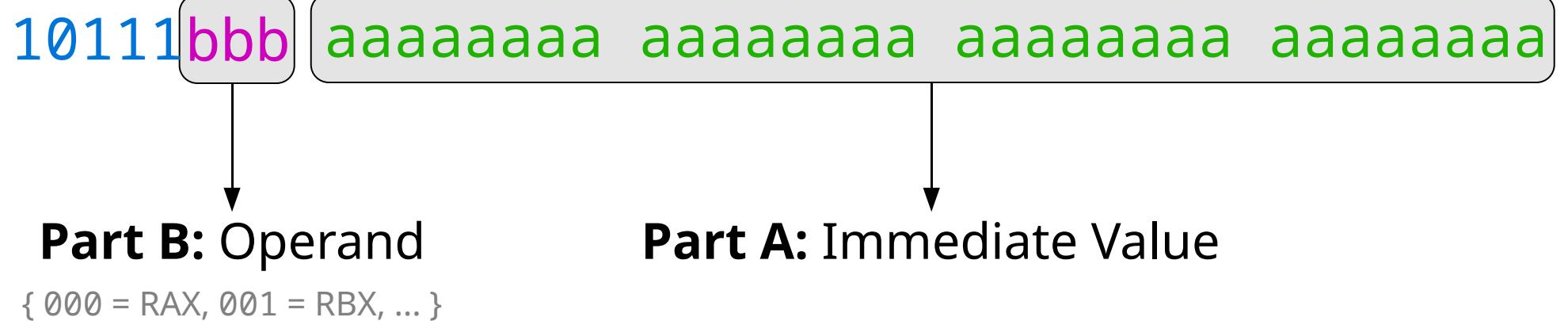
Dealing with long instructions



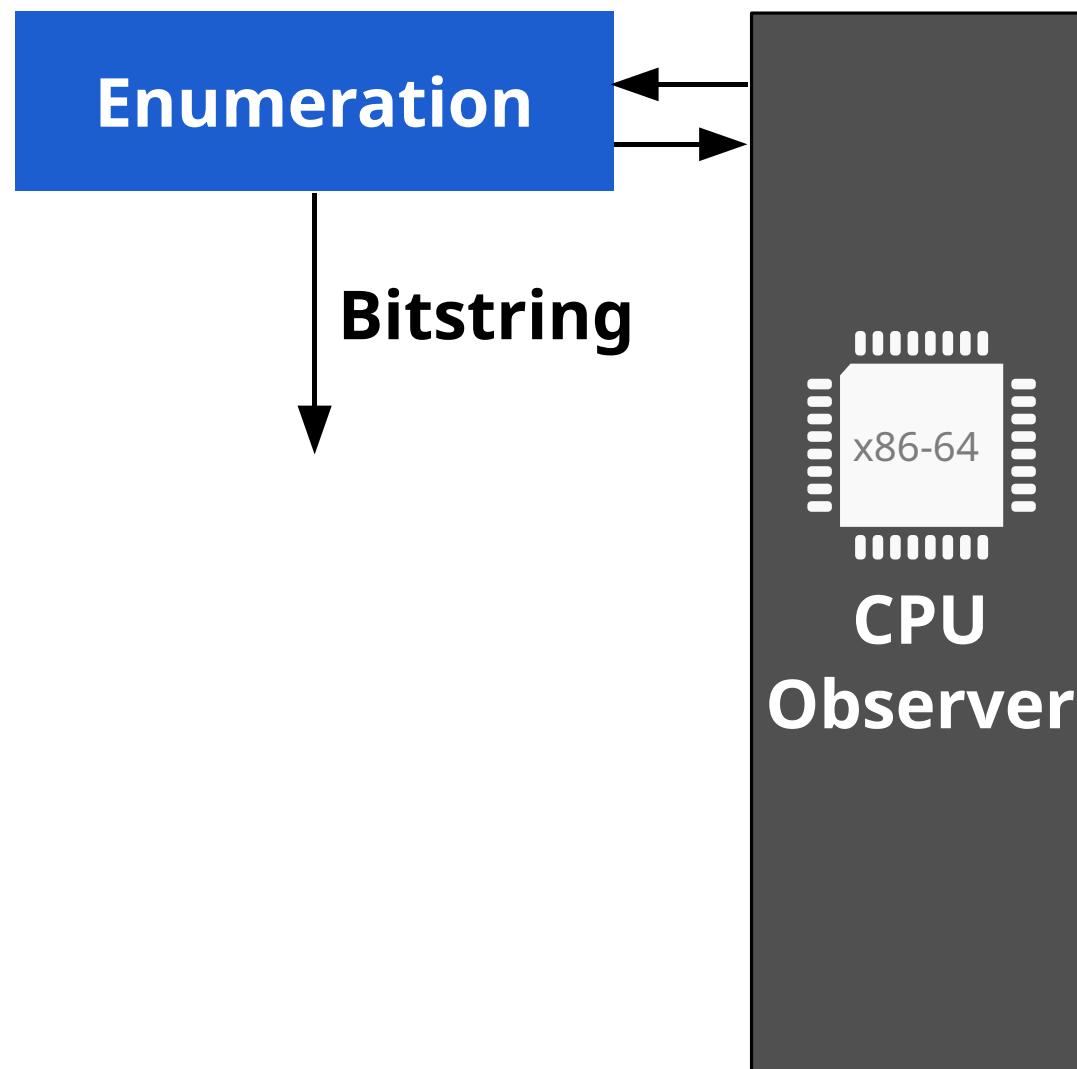
Encodings



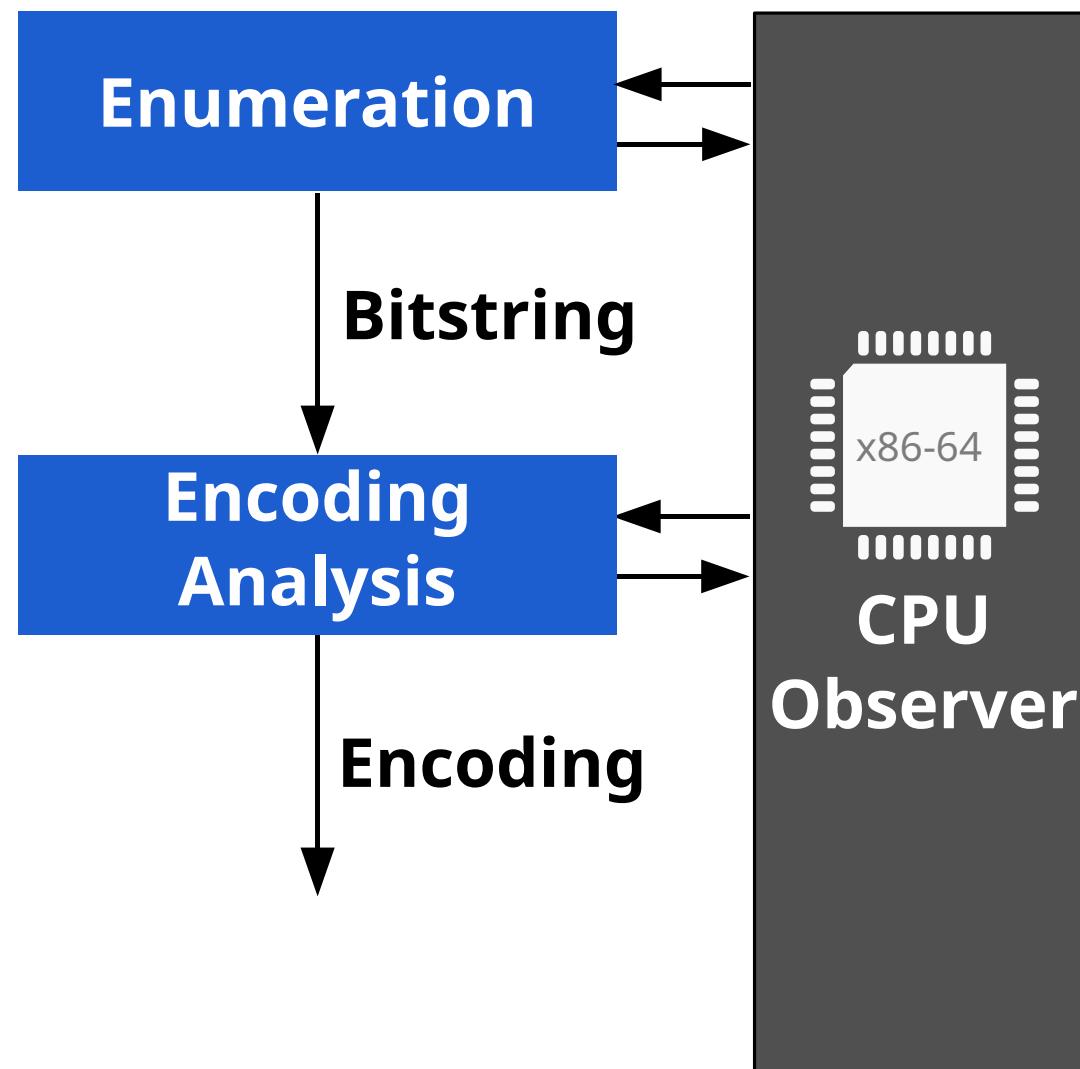
Encodings



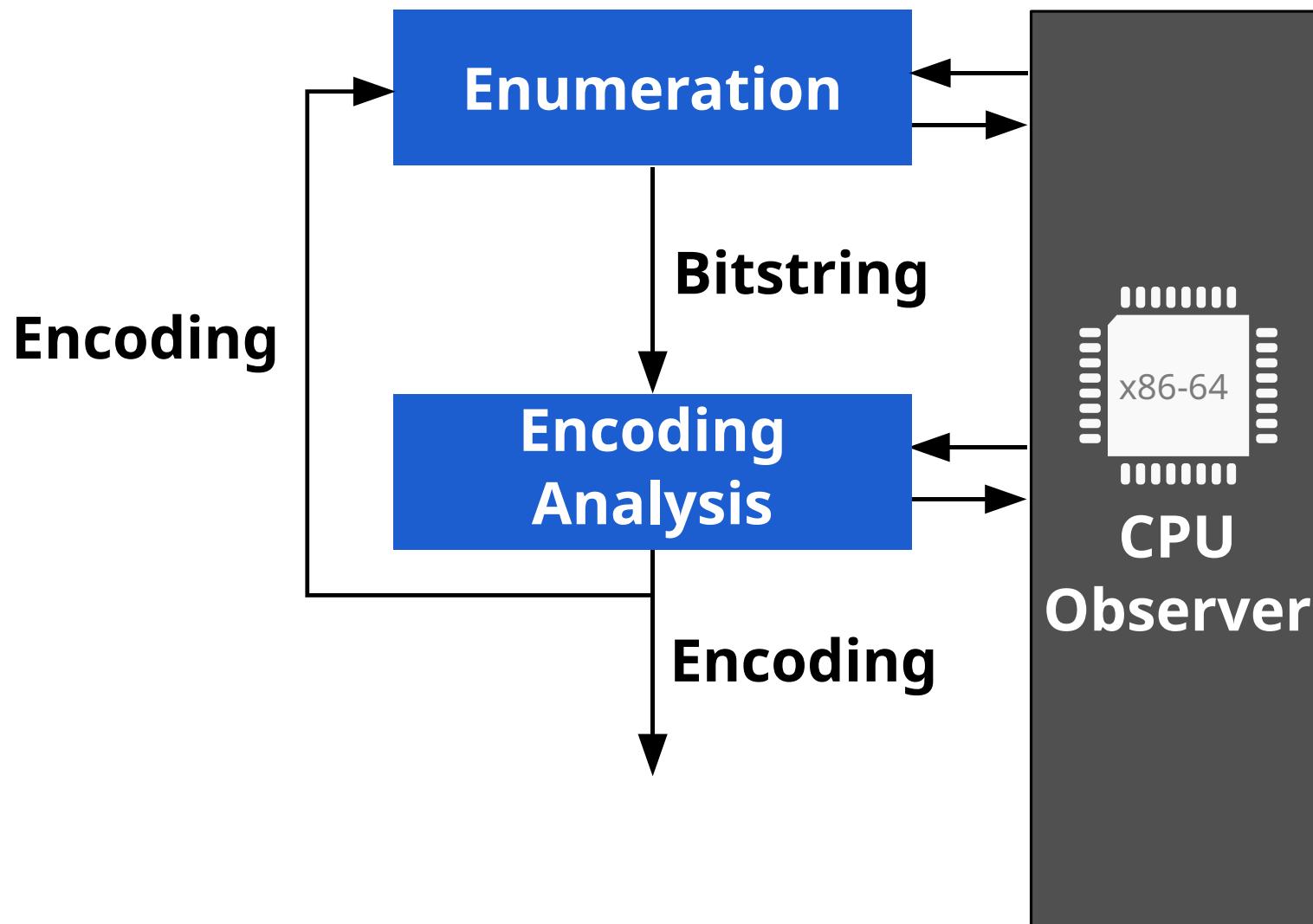
Analysis Overview



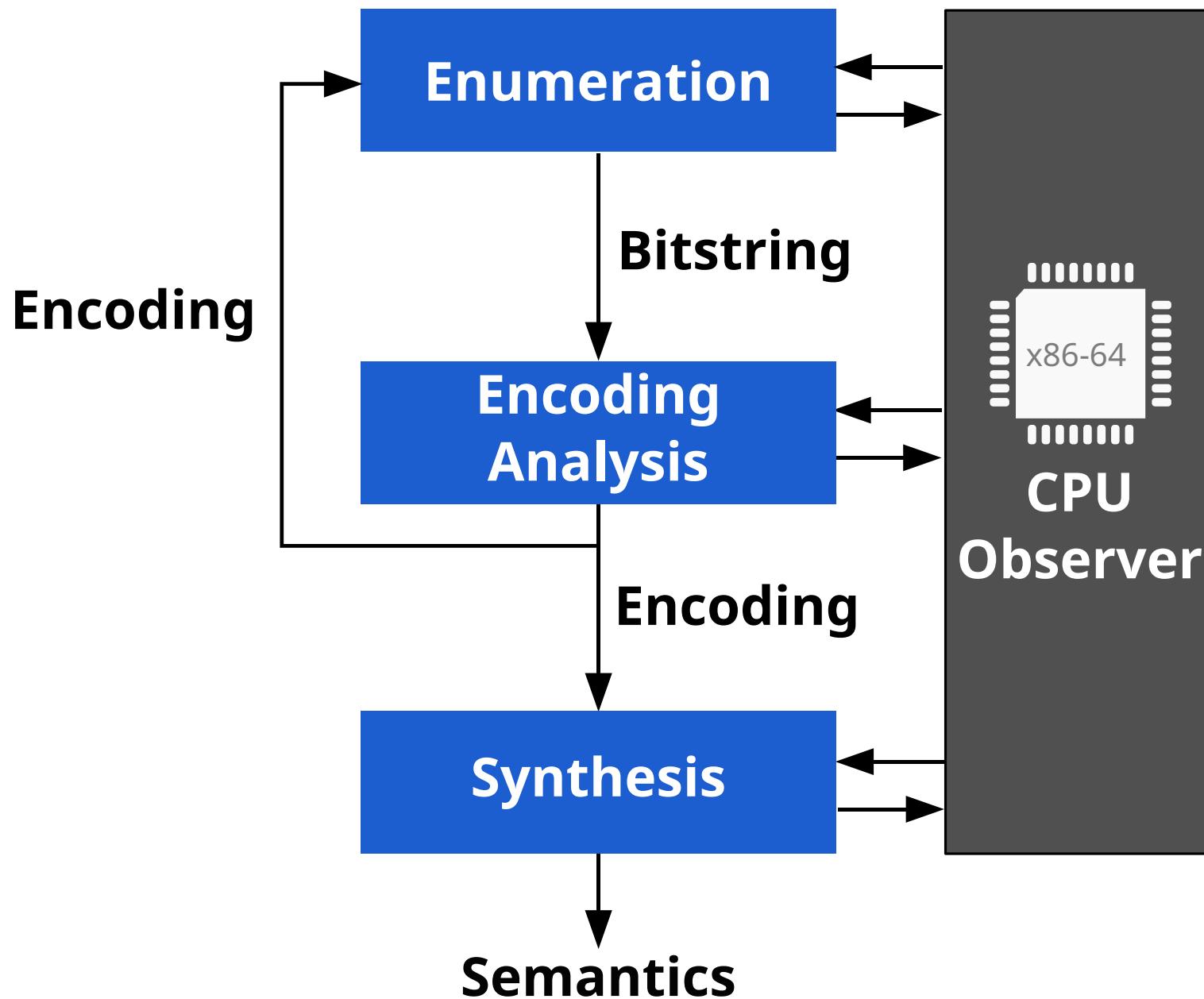
Analysis Overview



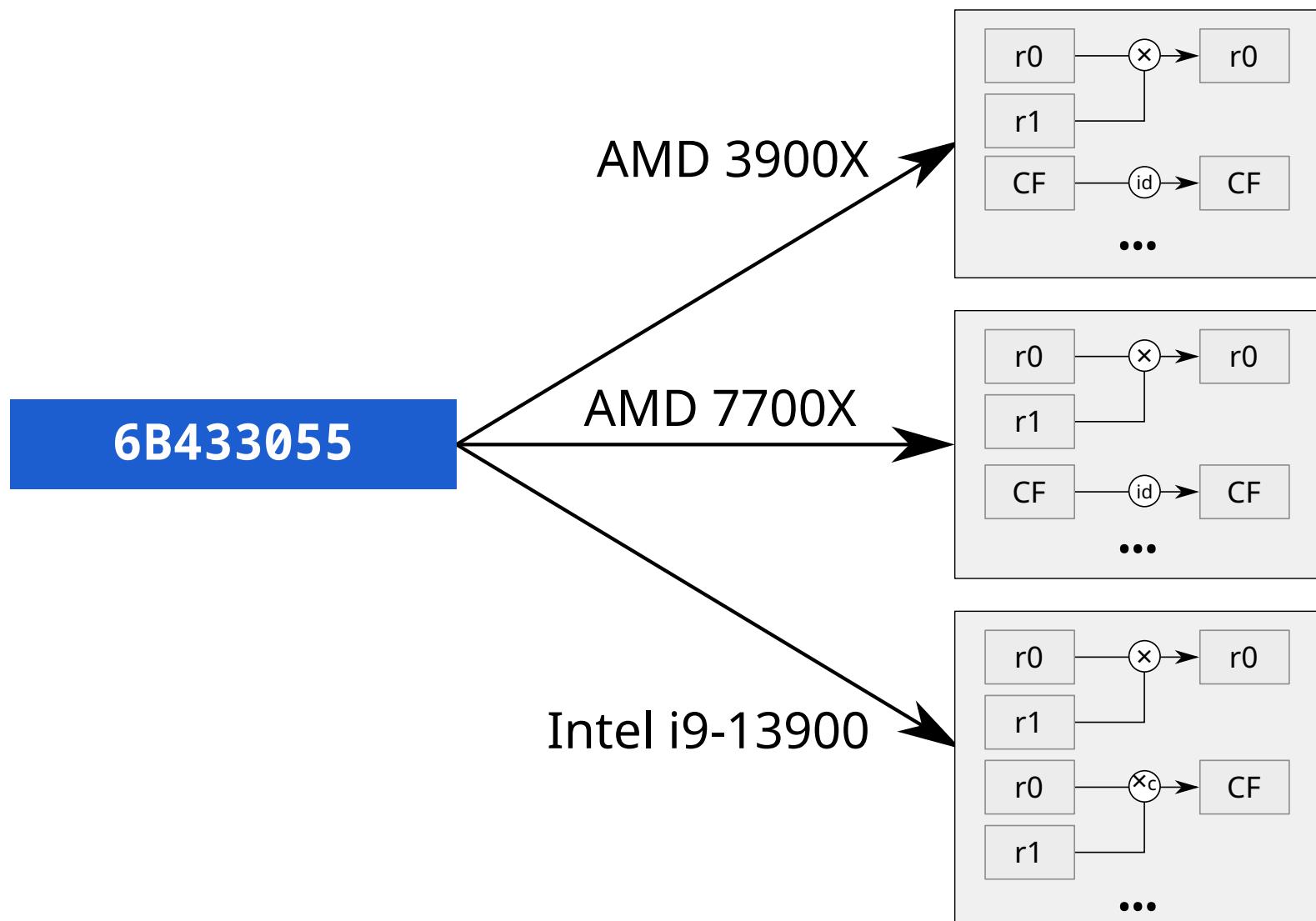
Analysis Overview



Analysis Overview



CPU-specific semantics



Analyzed CPUs

Architecture	A_0	A_1	A_2	A_3	A_4
CPU	AMD 3900X	AMD 7700X	Intel i9 13900 (p)	Intel i9 13900 (e)	Intel Xeon 4110
Cores	24T/12C	16T/8C	16T/8C	16T/16C	2x 16T/8C
Time taken	7 weeks	14 weeks	21 weeks	13 weeks	14 weeks
Encodings	118 025	118 019	117 605	118 135	117 229
Undocumented Encodings ¹	2693	3339	0	0	0
UB synthesized	91%	90%	91%	90%	90%

¹ copies of the VPERMQ* instructions with VEX.W = 0.

Analyzed CPUs

Architecture	A_0	A_1	A_2	A_3	A_4
CPU	AMD 3900X	AMD 7700X	Intel i9 13900 (p)	Intel i9 13900 (e)	Intel Xeon 4110
Cores	24T/12C	16T/8C	16T/8C	16T/16C	2x 16T/8C
Time taken	7 weeks	14 weeks	21 weeks	13 weeks	14 weeks
Encodings	118 025	118 019	117 605	118 135	117 229
Undocumented Encodings ¹	2693	3339	0	0	0
UB synthesized	91%	90%	91%	90%	90%

¹ copies of the VPERMQ* instructions with VEX.W = 0.

Analyzed CPUs

Architecture	A_0	A_1	A_2	A_3	A_4
CPU	AMD 3900X	AMD 7700X	Intel i9 13900 (p)	Intel i9 13900 (e)	Intel Xeon 4110
Cores	24T/12C	16T/8C	16T/8C	16T/16C	2x 16T/8C
Time taken	7 weeks	14 weeks	21 weeks	13 weeks	14 weeks
Encodings	118 025	118 019	117 605	118 135	117 229
Undocumented Encodings ¹	2693	3339	0	0	0
UB synthesized	91%	90%	91%	90%	90%

¹ copies of the VPERMQ* instructions with VEX.W = 0.

Analyzed CPUs

Architecture	A_0	A_1	A_2	A_3	A_4
CPU	AMD 3900X	AMD 7700X	Intel i9 13900 (p)	Intel i9 13900 (e)	Intel Xeon 4110
Cores	24T/12C	16T/8C	16T/8C	16T/16C	2x 16T/8C
Time taken	7 weeks	14 weeks	21 weeks	13 weeks	14 weeks
Encodings	118 025	118 019	117 605	118 135	117 229
Undocumented Encodings ¹	2693	3339	0	0	0
UB synthesized	91%	90%	91%	90%	90%

¹ copies of the VPERMQ* instructions with VEX.W = 0.

Comparing x86-64 implementations

Group	# Encodings	A_0	A_1	A_2	A_3	A_4
Group 0	95170	■	■	■	■	■
Group 1	4777	■	■	●	▲	●
Group 2	2571	■	■			
Group 3	1602	■	■	●	■	●
Group 4	604	■	●	▲	▼	○
Group 5	581	■	■	●	●	●
Group 6	101	■	■	■	■	
Group 7	40		■	■	■	
Group 8	29			■	■	
Group 9	24					■
Group 10	16	■	■	■	■	●
...

Comparing x86-64 implementations

Group	# Encodings	A_0	A_1	A_2	A_3	A_4
Group 0	95170	■	■	■	■	■
Group 1	4777	■	■	●	▲	●
Group 2	2571	■	■			
Group 3	1602	■	■	●	■	●
Group 4	604	■	●	▲	▼	○
Group 5	581	■	■	●	●	●
Group 6	101	■	■	■	■	
Group 7	40		■	■	■	
Group 8	29			■	■	
Group 9	24					■
Group 10	16	■	■	■	■	●
...

Comparing x86-64 implementations

Group	# Encodings	A_0	A_1	A_2	A_3	A_4
Group 0	95170	■	■	■	■	■
Group 1	4777	■	■	●	▲	●
Group 2	2571	■	■			
Group 3	1602	■	■	●	■	●
Group 4	604	■	●	▲	▼	○
Group 5	581	■	■	●	●	●
Group 6	101	■	■	■	■	
Group 7	40		■	■	■	
Group 8	29			■	■	
Group 9	24				■	
Group 10	16	■	■	■	■	●
...

Comparing x86-64 implementations

Group	# Encodings	A_0	A_1	A_2	A_3	A_4
Group 0	95170	■	■	■	■	■
Group 1	4777	■	■	●	▲	●
Group 2	2571	■	■			
Group 3	1602	■	■	●	■	●
Group 4	604	■	●	▲	▼	○
Group 5	581	■	■	●	●	●
Group 6	101	■	■	■	■	
Group 7	40		■	■	■	
Group 8	29			■	■	
Group 9	24					■
Group 10	16	■	■	■	■	●
...

Comparing x86-64 implementations

Group	# Encodings	A_0	A_1	A_2	A_3	A_4
Group 0	95170	■	■	■	■	■
Group 1	4777	■	■	●	▲	●
Group 2	2571	■	■			
Group 3	1602	■	■	●	■	●
Group 4	604	■	●	▲	▼	○
Group 5	581	■	■	●	●	●
Group 6	101	■	■	■	■	
Group 7	40		■	■	■	
Group 8	29			■	■	
Group 9	24					■
Group 10	16	■	■	■	■	●
...

Fingerprinting CPUs

```
xor rdx, rdx
xor rax, rax
add dl, 1
mul eax
setp r13b
setz r14b
lea rdi, [rip + 2f]
xor r15b, r15b
vpclmulhq1qdq ymm0, ymm1, ymm2
mov r15b, 1
```

2:

...

Fingerprinting CPUs

```
xor rdx, rdx
xor rax, rax
add dl, 1
mul eax
setp r13b
setz r14b
lea rdi, [rip + 2f]
xor r15b, r15b
vpc1mulhq1qdq ymm0, ymm1, ymm2
mov r15b, 1
```

2:

...

Semantics available on: explore.liblisa.nl

The screenshot shows the explore.liblisa.nl web application interface. At the top, it displays the URL <https://explore.liblisa.nl/instruction/00D3>. Below the URL, a banner indicates the processor architecture: 3900X | 7700X | i9-13900-p | i9-13900-e | Intel-Xeon-Silver-4110.

The main content area shows the assembly instruction `00000000 110bb0aa`. Below this, the "Parts" section displays two memory access diagrams:

- a**: Shows four memory locations: Rax (00), Rcx (01), Rdx (10), and Rbx (11). The Rdx and Rbx boxes are highlighted in red.
- b**: Shows four memory locations: Rax (00), Rcx (01), Rdx (10), and Rbx (11). The Rdx box is highlighted in green.

The "Memory accesses" section shows the address calculation: `Addr(m0) = Rip0:7` (2 bytes).

The "Output Dataflows SMTLib" section provides the semantic equations for the instruction:

- `Rip0:7 = 0x2 + sbe(Rip0:7)` (64 bit)
- `BL0 = sbe(DL0) + sbe(BL0)` (8 bit)
- `CF0 = Select[8:9](Crop[8](sbe(DL0)) + Crop[8](sbe(BL0))))` (1 bit)

Future work

Future work: validation of existing semantics

Instruction	libLISA	Das-gupta et al.	BAP	Ghidra	IDA (?)	...
0000000 110bb0aa	(assert (= ((_ extract 7 0) out_rRAX) (((_ extract 7 0) (_ extract 7 0)) (_ zero_extend 120) (_ extract 7 0) (bvadd (((_ sign_extend 120) (_ extract 7 0) (_ zero_extend 120) (_ extract 7 0) rRCX)))) (_ zero_extend 120) (_ extract 7 0) (_ zero_extend 120) (_ extract 7 0) rRAX)))))))		(assert (= rRAX (concat (((_ extract 63 8) rRAX) (((_ extract 7 0) (bvadd (concat #b0 (((_ extract 7 0) rRCX)) (concat #b0 (((_ extract 7 0) rRAX))))))))
10111011 11bbbaaa

Validation of BAP

Comparison report

Summary

This report describes the comparison results against [libLISA](#). Out of the **494** variants that were analyzed, errors were found in **211** variants:

ADC32mr	ADC32rm	ADC64mi8	ADC64rr_REV	BSR32rm	BSR32rr	BT32mr	BT64mr	BTC32mr	BTC64mr	BTR32mr
BTR64mr	BTS32mr	BTS64mr	CMPXCHG16B	FNSTCW16m	IMUL32m	IMUL32r	IMUL8m	IMUL8r	MMX_MOVD64from64rm	
MMX_MOVD64mr	MMX_MOVD64rm	MMX_MOVD64to64rm	MMX_MOVQ64mr	MMX_MOVQ64rm	MMX_MOVQ64rr_REV	MMX_PADD8irm				
MMX_PADDDirm	MMX_PADDWirm	MMX_PADDWrr	MMX_PALIGNRrmi	MMX_PALIGNRrrri	MMX_PANDNirm	MMX_PANDirm				
MMX_PAVGBirm	MMX_PAVGBirr	MMX_PAVGWirm	MMX_PCMPEQBirm	MMX_PCMPEQDirm	MMX_PCMPEQWirm	MMX_PCMPGTBirm				
MMX_PCMPGTDirm	MMX_PCMPGTDrr	MMX_PMAXSwirm	MMX_PMAXUBirm	MMX_PMINSWirm	MMX_PMINSWrr	MMX_PMINUBirm				
MMX_PORirm	MMX_PSLLDri	MMX_PSLLDrm	MMX_PSLLQri	MMX_PSLLQrm	MMX_PSLLQrr	MMX_PSLLWri	MMX_PSLLWrm			
MMX_PSLLWrr	MMX_PSRADrm	MMX_PSRLDri	MMX_PSRLDrm	MMX_PSRLQri	MMX_PSRLQrm	MMX_PSRLWrm	MMX_PSRLWrr			
MMX_PSUBBirm	MMX_PSUBQirm	MMX_PSUBWirm	MMX_PUNPCKHBWirm	MMX_PUNPCKHBWrr	MMX_PUNPCKHDQirm					
MMX_PUNPCKHWDirm	MMX_PUNPCKLBWirm	MMX_PUNPKLDDQirm	MMX_PUNPKLWDirm	MMX_PXORirm	MOV16ms	MOV32rs				
NEG8m	NOT32m	POP64rmm	ROL32m1	ROL32mCL	ROL32mi	ROL32rCL	ROL64m1	ROL64mi	ROL8mCL	ROL8mi
ROR32m1	ROR32mCL	ROR32mi	ROR64mi	ROR8m1	ROR8mi	ROR8r1	ROR8ri	SAR64rCL	SAR8mi	SAR8rCL
SBB32rm	SHL8mi	SHLD32mri8	SHLD32rrri8	SHLD64mri8	SHLD64rrri8	SHR8mi	SHRD32mri8	SHRD32rrri8		
SHRD64mri8	SHRD64rrri8	VMOVDDUPrm	VMOVSHDUPYrm	VMOVSHDUPrm	VMOVSLDUPYrm	VMOVSLDUPrm	VPADD8Yrm			
VPADD8rm	VPADD8Yrm	VPADDQYrm	VPADDQrm	VPADDQrr	VPADDWYrm	VPADDWrm	VPALIGNRrri	VPALIGNRrrri		
VPAVGBrm	VPAVGBrr	VPAVGWYrm	VPMAXSWYrm	VPMAXUBYrm	VPMAXUBYrr	VPMAXUBrm	VPMAXUDrm	VPMAXUDrr		
VPMAXUWYrm	VPMINSDYrm	VPMINSWYrm	VPMINSWYrr	VPMINSWrm	VPMINUBYrr	VPMINUBrm	VPMINUBrr	VPMINUDYrm		
VPMINUDrm	VPMINUWYrm	VPMINUWrm	VPSHUFDYmi	VPSHUFDYri	VPSHUFDmi	VPSHUFDri	VPSLLDQYri	VPSLLDYri		
VPSLLDri	VPSLLDrm	VPSLLQYri	VPSLLQYrm	VPSLLQri	VPSLLQrm	VPSLLWYrm	VPSLLWrm	VPSRADYrm	VPSRADrm	
VPSRADrr	VPSRAWYrm	VPSRAWrm	VPSRLDQYri	VPSRLDYri	VPSRLDYrm	VPSRLDYrr	VPSRLDri	VPSRLDrm	VPSRLQYri	
VPSRLQYrm	VPSRLQri	VPSRLQrm	VPSRLQrr	VPSRLWYrm	VPSRLWri	VPSRLWrm	VPSUBBYrm	VPSUBBYrr	VPSUBBrm	
VPSUBDrm	VPSUBQYrm	VPSUBQrm	VPSUBWYrm	VPUKPCKHBWrm	VPUKPCKHDQYrm	VPUKPCKHDQrr	VPUKPCKHQDQYrm			
VPUKPCKHQDQrm	VPUKPCKHWYrm	VPUKPCKHWDrm	VPUKPCKHWDrri	VPUKPCKLBWYrm	VPUKPCKLBWrm	VPUKPCKLDQYrm				
VPUKPCKLDQrm	VPUKPCKLDQrr	VPUKPCKLDQYrm	VPUKPCKLDYrm	VPUKPCKLWDrm	VXORPDrm	VXORPSYrm	VXORPSrm			
XADD32rm	XADD64rm	XCHG32rm								

Comparison report

Summary

T

- Out of the **494** variants that were analyzed, errors were found in **211** variants:

MMX_MOVD64mr	MMX_MOVD64rm	MMX_MOVD64to64rm	MMX_MOVQ64mr	MMX_MOVQ64rm	MMX_MOVQ64rr_REV	MMX_PADDBirm
MMX_PADDDirm	MMX_PADDWirm	MMX_PADDWrr	MMX_PALIGNRrmi	MMX_PALIGNRrrri	MMX_PANDNirm	MMX_PANDirm
MMX_PAVGBirm	MMX_PAVGBirr	MMX_PAVGWirm	MMX_PCMPEQBirm	MMX_PCMPEQDirm	MMX_PCMPEQWirm	MMX_PCMPGTBirm
MMX_PCMPGTDirm	MMX_PCMPGTDrr	MMX_PMAXSWirm	MMX_PMAXUBirm	MMX_PMINSWirm	MMX_PMINSWrr	MMX_PMINUBirm
MMX_PORirm	MMX_PSLLDri	MMX_PSLLDrm	MMX_PSLLQri	MMX_PSLLQrm	MMX_PSLLQrr	MMX_PSLLWri
MMX_PSLLWrr	MMX_PSRADrm	MMX_PSRLDri	MMX_PSRLDrm	MMX_PSRLQri	MMX_PSRLQrm	MMX_PSRLWrm
MMX_PSUBBirm	MMX_PSUBQirm	MMX_PSUBWirm	MMX_PUNPCKHBWirm	MMX_PUNPCKHBWrr	MMX_PUNPCKHDQirm	
MMX_PUNPCKHWDirm	MMX_PUNPCKLBWirm	MMX_PUNPCKLDQirm	MMX_PUNPCKLWDirm	MMX_PXORirm	MOV16ms	MOV32rs
NEG8m	NOT32m	POP64rmm	ROL32m1	ROL32mCL	ROL32mi	ROL32rCL
ROL64m1	ROL64mi	ROL64mCL	ROL64mi	ROL64m1	ROL8mCL	ROL8mi
ROL8ri	ROR32m1	ROR32mCL	ROR32mi	ROR64mi	ROR8m1	ROR8mi
ROR8ri	SAR64rCL	SAR8mi	SAR8rCL	SAR8m1	SAR8ri	SBB32mr
SBB32rm	SHL8mi	SHLD32mri8	SHLD32rrri8	SHLD64mri8	SHLD64rrri8	SHR8mi
SHRD32mri8	SHRD32rrri8	VMOVDDUPrm	VMOVSHDUPYrm	VMOVSHDUPrm	VMOVSLDUPYrm	VMOVSLDUPrm
VPADDBrm	VPADDYrm	VPADDQYrm	VPADDQrm	VPADDQrr	VPADDWYrm	VPADDWrm
VPALIGNRrri	VPALIGNRYrrri	VPAVGBrm	VPAVGBrr	VPAVGWYrm	VPMAXSWYrm	VPMAXUBYrm
VPMAXUBYrr	VPMAXUBrm	VPMAXUDrm	VPMAXUDrr	VPMAXUWYrm	VPMINSDYrm	VPMINSWYrm
VPMINSDYrr	VPMINSWrm	VPMINSWYrr	VPMINSWYrm	VPMINUBYrr	VPMINUBYrm	VPMINUBrm
VPMINUBrr	VPMINUDrm	VPMINUWYrm	VPMINUWrm	VPMINUWYri	VPMINUWYrm	VPMINUWrr
VPMINUWrr	VPMINUWYrm	VPMINUWrm	VPSHUFDYmi	VPSHUFDYri	VPSHUFDmi	VPSHUFDri
VPSHUFDri	VPSHUFDrr	VPSHUFDYri	VPSHUFDYrm	VPSHUFDYrr	VPSHUFDYrm	VPSHUFDYrr
VPSLLDri	VPSLLDrm	VPSLLQYri	VPSLLQYrm	VPSLLQri	VPSLLQrm	VPSLLWYrm
VPSLLWrm	VPSRADYrm	VPSRADrr	VPSRAWYrm	VPSRAWrm	VPSRLDYri	VPSRLDYrm
VPSRLDYrm	VPSRLQri	VPSRLQrm	VPSRLQrr	VPSRLWYrm	VPSRLWri	VPSRLWrm
VPSRLWrm	VPSUBBYrm	VPSUBBYrr	VPSUBBYrr	VPSUBBYrm	VPSUBBYrr	VPSUBBrm
VPSUBDrm	VPSUBQYrm	VPSUBQrm	VPSUBWYrm	VPUNPCKHBWrm	VPUNPCKHDQYrm	VPUNPCKHDQrm
VPUNPCKHDQrm	VPUNPCKHWDrm	VPUNPCKHWDrri	VPUNPCKHWYrm	VPUNPCKLBWrm	VPUNPCKLBWrr	VPUNPCKLDQYrm
VPUNPCKLDQrm	VPUNPCKLDQrr	VPUNPCKLQDQYrm	VPUNPCKLWDYrm	VPUNPCKLWDrm	VXORPDrm	VXORPSYrm
VXORPSYrm	VXORPSrm	XADD32rm	XADD64rm	XCHG32rm		

Variant VPADDQYrm

Instruction C55DD45941 decoded as VPADDQYrm(YMM11, YMM4, RCX, 0x1, Nil, 0x41, Nil)

- ▼ The location **Xmm4** should not be modified.

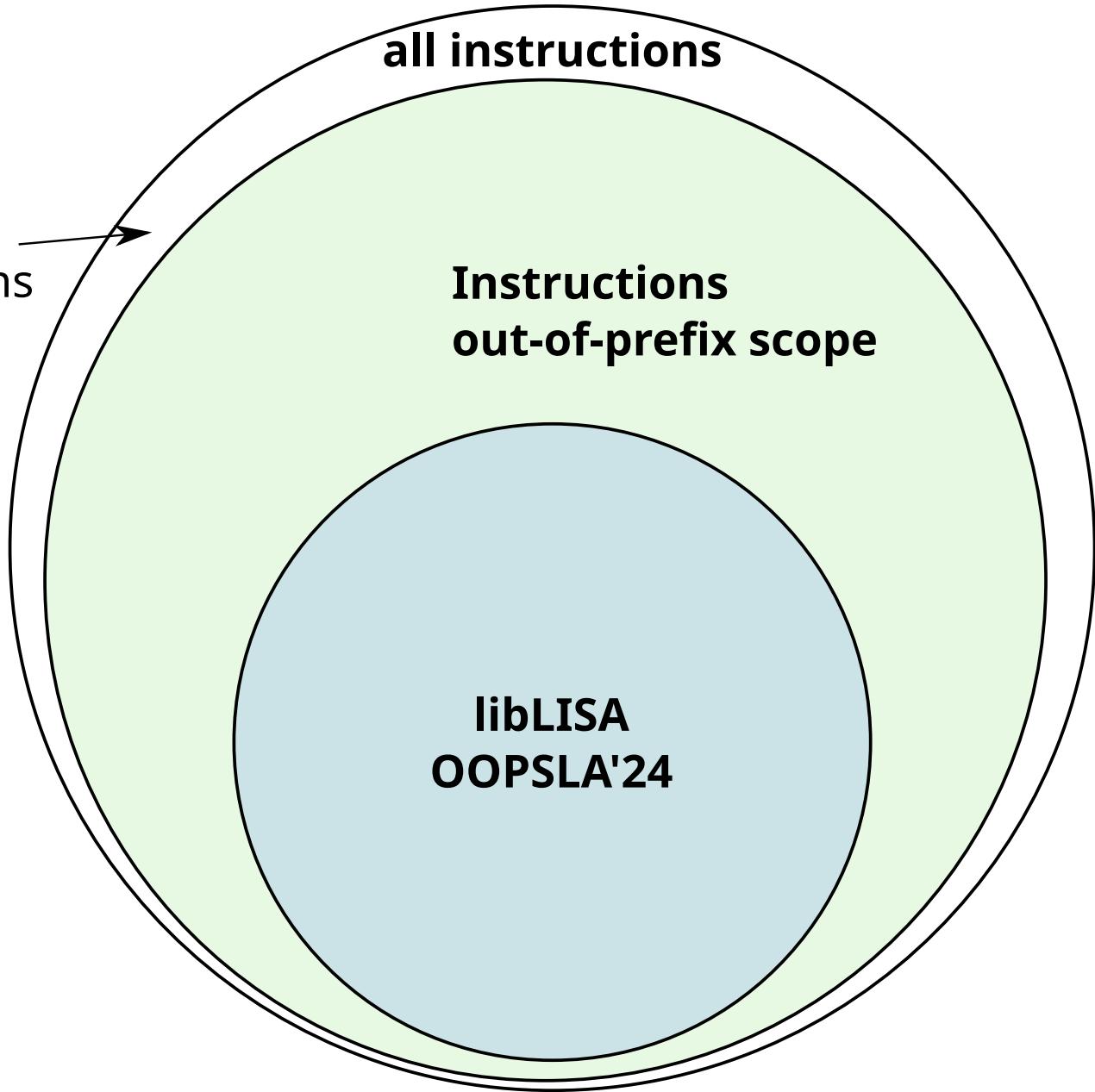
- ▶ The location **Xmm11** should be updated.

Instruction C4E13DD4ACB3F44309F2 decoded as VPADDQYrm(YMM5, YMM8, RBX, 0x4, RSI, -0xdff6bc0c, Nil)

- ▶ The location **Xmm5** should be updated.
 - ▶ The location **Xmm8** should not be modified.

Future work: prefixes

Privileged instructions
Hardware instructions
Segment selector instructions
(rdtsc, xsave, out, lss, lsl, ...)



Future work: prefixes

0303 = add eax, dword ptr [rbx]

480303 = add **rax**, qword ptr [rbx]

660303 = add **ax**, word ptr [rbx]

Future work: prefixes

0303 = add eax, dword ptr [rbx]

480303 = add rax, qword ptr [rbx]

660303 = add ax, word ptr [rbx]

48486648F2F348662E3E4848660303 = add ax, word ptr [rbx]

Summary



- automatically infers semantics
- analyzed 5 architectures
- semantics and code are open source

