# Prefetching in a Blockchain Client using static analysis and speculative execution of Smart Contracts

Diploma Thesis
Panagiotis Gkonis

# Ethereum clients

- Ethereum network participants run a client

- It downloads the entire blockchain and runs all contained transaction

- A very time consuming process: runs for hours or even days on modern hardware

- Final disk usage: a few TB

# World State

- It's a set of K-V pairs, identical on all clients, modified by Transactions

- Contains accounts:
  Key=address → Value=account

- and for each contract, its storage:
  Key=slot → Value=content

- Structured as a modified Merkle Patricia Trie (MPT)

# Merkle Patricia Trie

- Functions as a KV store

- Tree-like data structure

- The key is the path from root to a given node

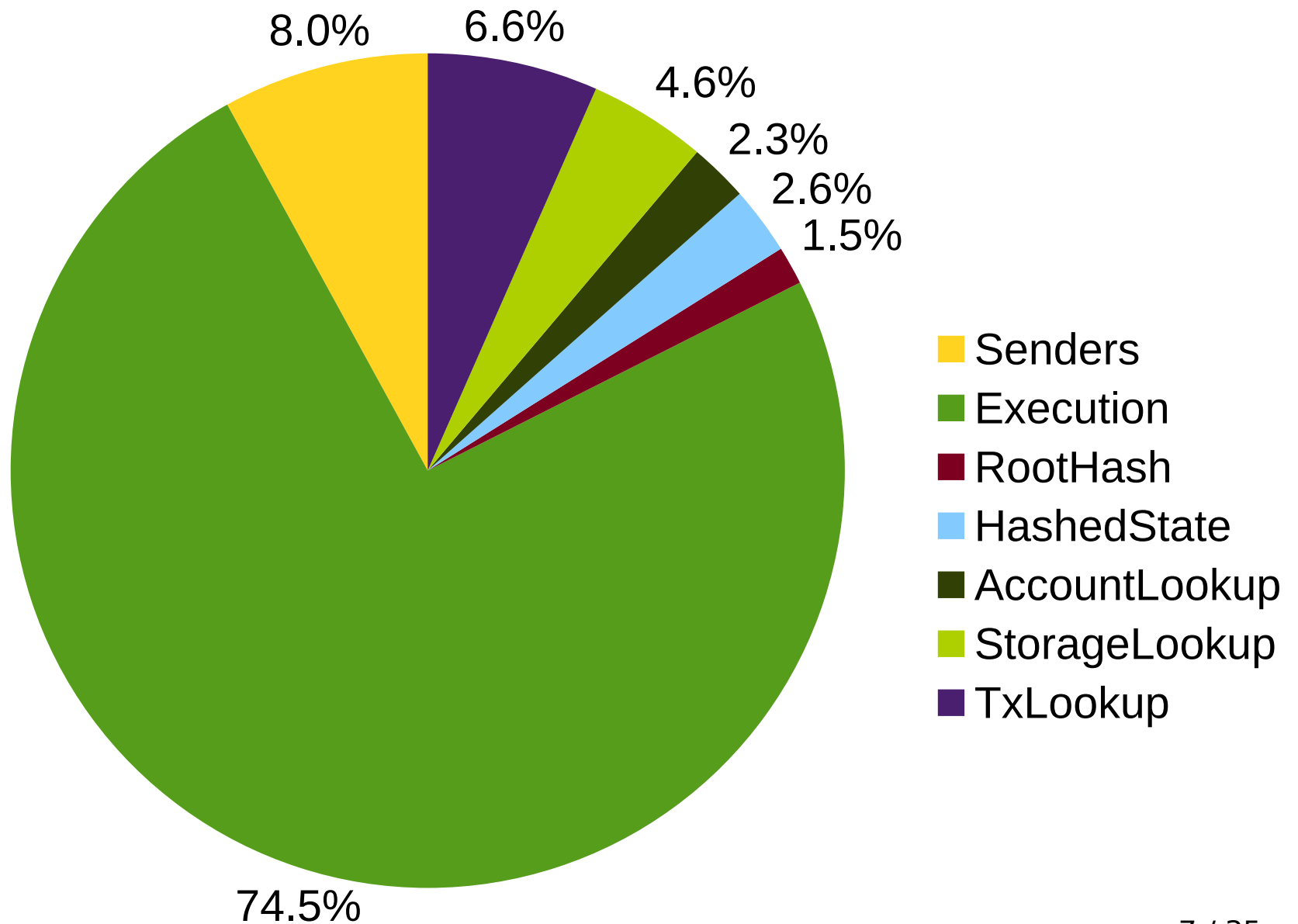- That node contains the corresponding value

# Go-ethereum

- Official client

- Internally implements a MPT

- which is stored in LevelDB

- Accessing historic state is trivial

- DB size of an archive node (2022): 9 TB

# Erigon

- Fork of go-ethereum

- Works in stages, e.g. download, execute, …

- Does not use a Trie, instead storing in a "flat" data structure

- Database: MDBX (a fork of LMDB), mmap-ed

- Considerable improvement in space and time

- DB size of an archive node (2022): 1.5 TB

# Stages (excl. network)



Senders 8.0%
TxLookup 6.6%
StorageLookup 4.6%
AccountLookup 2.3%
HashedState 2.6%
RootHash 1.5%
Execution 74.5%

Legend:
- Senders
- Execution
- RootHash
- HashedState
- AccountLookup
- StorageLookup
- TxLookup

# Execution Stage

- Our main focus: it executes blocks, transactions and smart contracts (SC)

- Sequential execution, due to unknown inter-transactional dependencies

- CPU and IO heavy

- Constant access to World State (DB), with time-consuming blocking reads

- If we used prefetching, the DB's pages would be cached in memory (FS cache), accelerating reads
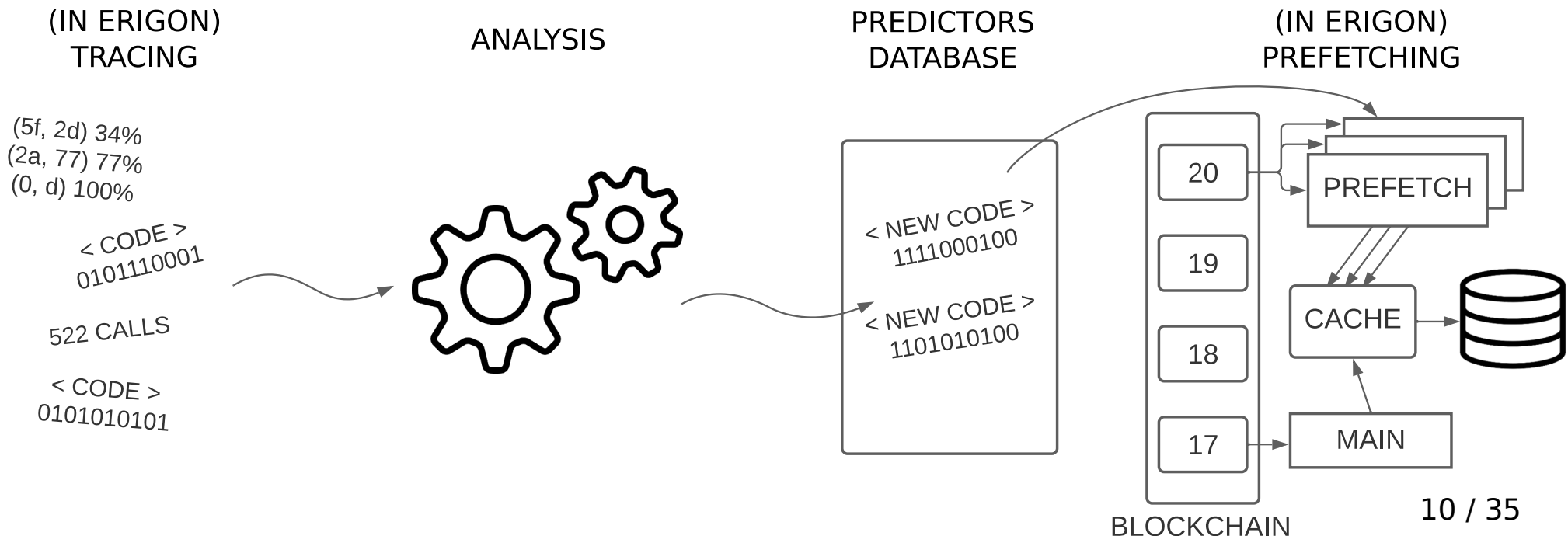
  → precisely the purpose of this thesis

# How much I/O;

- Test run of 30K blocks, measuring percentage of time spent in DB queries

➔ With a fast NVME drive:   22 %

➔ With a slower sata SSD:   35 %

- There is an oportunity for considerable improvement

*(no test was performed with a hard drive, as it would be 1-2 orders of magnitude slower)*

# System overview

- Tracing: collecting metrics and SC code

- Anslysis of collected SC and synthesis of new micro-programs (predictors)

- Prefetching and speculative execution of predictors



(IN ERIGON)
TRACING

ANALYSIS

PREDICTORS
DATABASE

(IN ERIGON)
PREFETCHING

(5f, 2d) 34%
(2a, 77) 77%
(0, d) 100%

< CODE >
0101110001

522 CALLS

< CODE >
0101010101

< NEW CODE >
1111000100

< NEW CODE >
1101010100
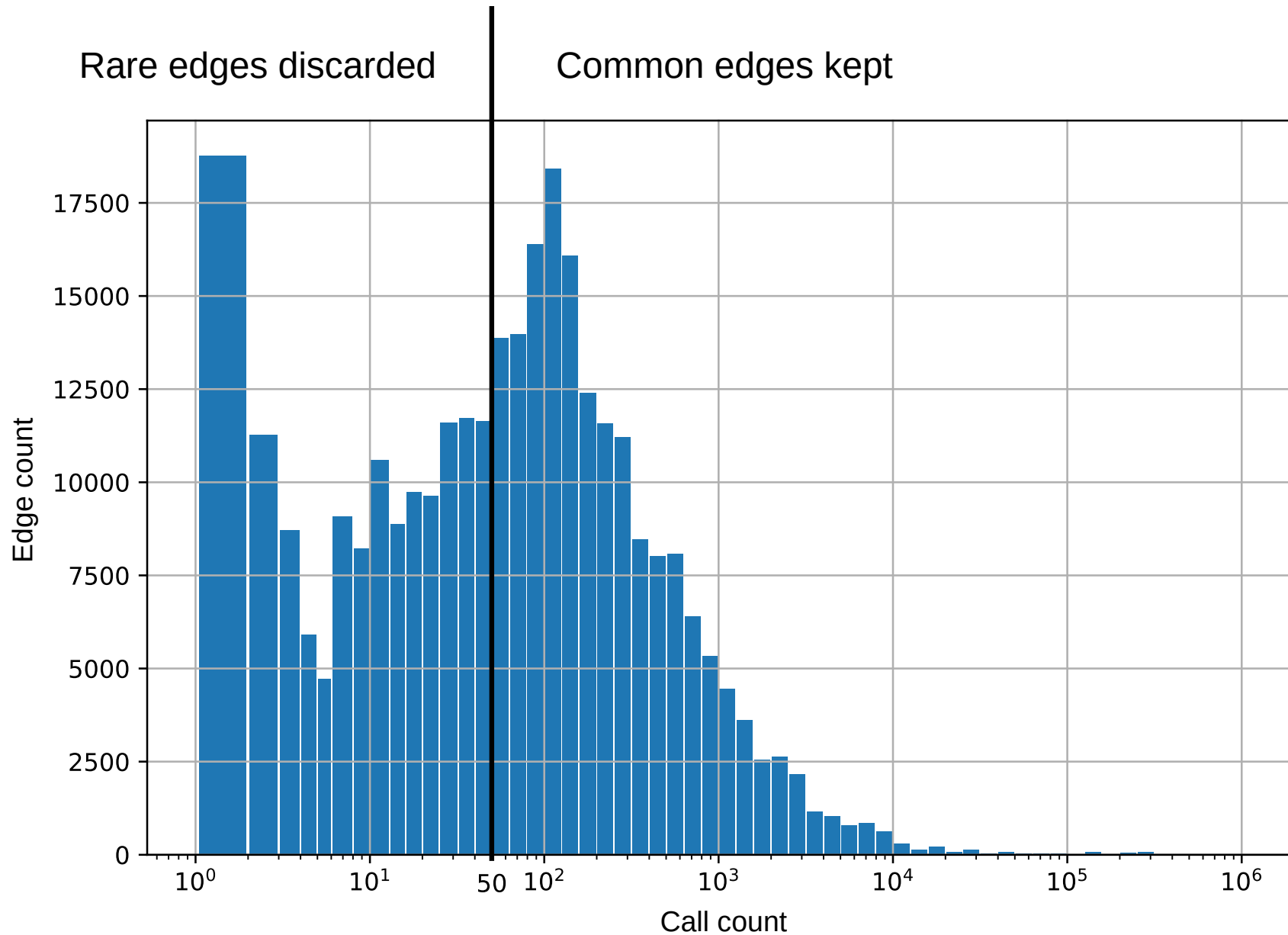
20

19

18

17

PREFETCH

CACHE

MAIN

BLOCKCHAIN

# Tracing

- During the execution stage*, the SC code is collected and the number of calls to each one

- Frequent SCs are saved for analysis

- Also collected, for each JUMP, the PC before and after, as well as the unique ID of the call being executed

- These JUMP address pairs (CFG edges) with "many" calls are given as hints to the analyzer

*here realized with a transaction sample, offline, though in a real implementation would be online

# Tracing

# Analysis

- Separate process, in Python

- Input is SC code, output is micro-programs (predictors) with similar code

- Predictors contain only "useful" instructions: storage access, calls, etc.

- Behavior can differ from the original code

- Opportunities for unsafe optimization, e.g. removing rare instructions

# Analysis 1/12

- Takes as input the code of a SC, which is in EVM bytecode

608060405260043610610056576000000000000000000000000000000000000000000000000000000000000000006000350416636ae17ab7811461005b578063771c0ad91461008d578063e80db5db146100ab575b600080fd5b34801561006757600080f...

# Analysis 2/12

- It disassembles

```
00000000: PUSH1 0x80
00000002: PUSH1 0x40
00000004: MSTORE
00000005: PUSH1 0x4
00000007: CALLDATASIZE
00000008: LT
00000009: PUSH2 0x56
0000000c: JUMPI
0000000d: PUSH4 0xffffffff
00000012: PUSH29 0x10000000000000...
00000030: PUSH1 0x0
00000032: CALLDATALOAD
00000033: DIV
00000034: AND
....
```

# Analysis 3/12

- Breaks it down into basic blocks

```
-----     BLOCK ~0 -----
00000000: PUSH1 0x80
00000002: PUSH1 0x40
00000004: MSTORE
00000005: PUSH1 0x4
00000007: CALLDATASIZE
00000008: LT
00000009: PUSH2 0x56
0000000c: JUMPI
```

```
-----     BLOCK ~d -----
0000000d: PUSH4 0xffffffff
00000012: PUSH29 0x100000000...
00000030: PUSH1 0x0
00000032: CALLDATALOAD
00000033: DIV
00000034: AND
00000035: PUSH4 0x6ae17ab7
0000003a: DUP2
0000003b: EQ
0000003c: PUSH2 0x5b
0000003f: JUMPI
```

```
-----     BLOCK ~40 -----
00000040: DUP1
00000041: PUSH4 0x771c0ad9
00000046: EQ
00000047: PUSH2 0x8d
0000004a: JUMPI
```

# Analysis 4/12

- Converts instructions to SSA form

```
-----      BLOCK ~0 -----
0x0:  .3 = PHI~0-MEM
0x0:  .0 = #80
0x2:  .1 = #40
0x4:  .2 = MSTORE(.3, .1, .0)
0x5:  .4 = #4
0x7:  .5 = CALLDATASIZE
0x8:  .6 = LT(.5, .4)
0x9:  .7 = #56
0xc:  .8 = JUMPI(.7, .6)
```

```
-----      BLOCK ~d -----
0xd:  .0 = #ffffffff
0x12: .1 = #10000...
0x30: .2 = #0
0x32: .3 = CALLDATALOAD(.2)
0x33: .4 = DIV(.3, .1)
0x34: .5 = AND(.4, .0)
0x35: .6 = #6ae17ab7
0x3b: .7 = EQ(.5, .6)
0x3c: .8 = #5b
0x3f: .9 = JUMPI(.8, .7)
```

```
-----      BLOCK ~40 -----
0x40: .0 = PHI~40[-1]
0x41: .1 = #771c0ad9
0x46: .2 = EQ(.1, .0)
0x47: .3 = #8d
0x4a: .4 = JUMPI(.3, .2)
```

# Analysis 5/12

- Connects blocks,
  initial CFG estimation

```
-----    BLOCK ~0 -----
0x0: .3 = PHI~0-MEM
0x0: .0 = #80
0x2: .1 = #40
0x4: .2 = MSTORE(.3, .1, .0)
0x5: .4 = #4
0x7: .5 = CALLDATASIZE
0x8: .6 = LT(.5, .4)
0x9: .7 = #56
0xc: .8 = JUMPI(.7, .6)
```

NT

```
-----    BLOCK ~d -----
0xd: .0 = #ffffffff
0x12: .1 = #100000...
0x30: .2 = #0
0x32: .3 = CALLDATALOAD(.2)
0x33: .4 = DIV(.3, .1)
0x34: .5 = AND(.4, .0)
0x35: .6 = #6ae17ab7
0x3b: .7 = EQ(.5, .6)
0x3c: .8 = #5b
0x3f: .9 = JUMPI(.8, .7)
```

NT                    T

```
-----    BLOCK ~40 -----
0x40: .0 = PHI~40[-1](~d.5)
0x41: .1 = #771c0ad9
0x46: .2 = EQ(.1, .0)
0x47: .3 = #8d
0x4a: .4 = JUMPI(.3, .2)
```

```
-----    BLOCK ~5b -----
0x5c: .0 = CALLVALUE
0x5e: .1 = ISZERO(.0)
0x5f: .2 = #67
0x62: .3 = JUMPI(.2, .1)
```

- Performs optimizations, using the worklist algorithm

```
-----     BLOCK ~0 -----
0x0:  .3 = PHI~0-MEM
0x0:  .0 = #80
0x2:  .1 = #40
0x4:  .2 = MSTORE(.3, .1, .0)
0x5:  .4 = #4
0x7:  .5 = CALLDATASIZE
0x8:  .6 = LT(.5, .4)
0x9:  .7 = #56
0xc:  .8 = JUMPI(.7, .6)
```

NT

```
-----     BLOCK ~d -----
0xd:  .10 = PHI~d-MEM(~0.2)
0xd:  .0 = #ffffffff
0x12: .1 = #10000...
0x30: .2 = #0
0x32: .3 = CALLDATALOAD(.2)
0x33: .4 = DIV(.3, .1)
0x34: .5 = AND(.4, .0)
0x35: .6 = #6ae17ab7
0x3b: .7 = EQ(.5, .6)
0x3c: .8 = #5b
0x3f: .9 = JUMPI(.8, .7)
```

```
-----     BLOCK ~8d -----
0x8e: .0 = CALLVALUE
0x90: .1 = ISZERO(.0)
0x91: .2 = #99
0x94: .3 = JUMPI(.2, .1)
```

NT

b_95

```
-----     BLOCK ~99 -----
0x99: .0 = PHI~99[-1]
0x9b: .1 = #79
0x9e: .2 = #...
0xa0: .3 = CALLDATALOAD(.2)
0xa1: .4 = #24
0xa3: .5 = CALLDATALOAD(.4)
0xa4: .6 = #44
0xa6: .7 = CALLDATALOAD(.6)
0xa7: .8 = #10f
0xaa: .9 = JUMP(.8)
```

NT                    T

```
-----     BLOCK ~40 -----
0x40: .0 = PHI~40[-1](~d.5)
0x41: .1 = #771c0ad9
0x46: .2 = EQ(.1, .0)
0x47: .3 = #8d
0x4a: .4 = JUMPI(.3, .2)
```

```
-----     BLOCK ~5b -----
0x5b: .4 = PHI~5b-MEM(~d.10)
0x5c: .0 = CALLVALUE
0x5e: .1 = ISZERO(.0)
0x5f: .2 = #67
0x62: .3 = JUMPI(.2, .1)
```

# Analysis 7/12

- Chooses instructions to be included in the predictor, e.g. SLOAD, CALL and dependencies

```
----- * BLOCK ~0 -----
*0x0:  .3 \ PHI~0-MEM
0x0:   .0 = #80
0x2:   .1 = #40
*0x4:  .2 \ MSTORE(.3, .1#40, .0#80)
0x5:   .4 = #4
0x7:   .5 = CALLDATASIZE
0x8:   .6 = LT(.5, .4#4)
0x9:   .7 = #56
0xc:   .8 \ JUMPI(.7#56, .6)
```

NT

```
----- * BLOCK ~d -----
*0xd:  .10 \ PHI~d-MEM(~0.2)
0xd:   .0 = #ffffffff
0x12:  .1 = #10000...
0x30:  .2 = #0
*0x32: .3 = CALLDATALOAD(.2#0)
*0x33: .4 = DIV(.3, .1#1000)
*0x34: .5 = AND(.4, .0#ffff)
0x35:  .6 = #6ae17ab7
*0x3b: .7 = EQ(.5, .6#6ae1)
0x3c:  .8 = #5b
*0x3f: .9 \ JUMPI(.8#5b, .7)
```

NT          T

```
-----   BLOCK ~40 -----
0x40:  .0 = PHI~40[-1](~d.5)
0x41:  .1 = #771c0ad9
0x46:  .2 = EQ(.1#771c, .0)
0x47:  .3 = #8d
0x4a:  .4 \ JUMPI(.3#8d, .2)
```

```
----- * BLOCK ~5b -----
*0x5b: .4 \ PHI~5b-MEM(~d.10)
*0x5c: .0 = CALLVALUE
*0x5e: .1 = ISZERO(.0)
0x5f:  .2 = #67
*0x62: .3 \ JUMPI(.2#67, .1)
```

# Analysis 8/12

- Enumerates values of chosen instructions

```
----- ON MAP -----
1 = #40
2 = #80
3 = V~0.2-MSTORE(v~0.3-PHIxb232-0B, #40, #80)-xad80-NV
4 = #10000...
5 = #0
6 = #6ae17ab7
7 = #5b
8 = #ffffffff
9 = V~d.3-CALLDATALOAD(#0)-x15b2
10 = V~d.4-DIV(v~d.3-CALLDATALOADx15b2, #10000...)-x4ea2
11 = V~d.5-AND(v~d.4-DIVx4ea2, #ffffffff)-x4954
12 = V~d.7-EQ(v~d.5-ANDx4954, #6ae17ab7)-x30c9
13 = V~d.9-JUMPI(#5b, v~d.7-EQx30c9)-x2f1e-NV
14 = #67
15 = V~5b.0-CALLVALUE()-x78d0
16 = V~5b.1-ISZERO(v~5b.0-CALLVALUEx78d0)-x8a44
17 = V~5b.3-JUMPI(#67, v~5b.1-ISZEROx8a44)-x9d52-NV
```
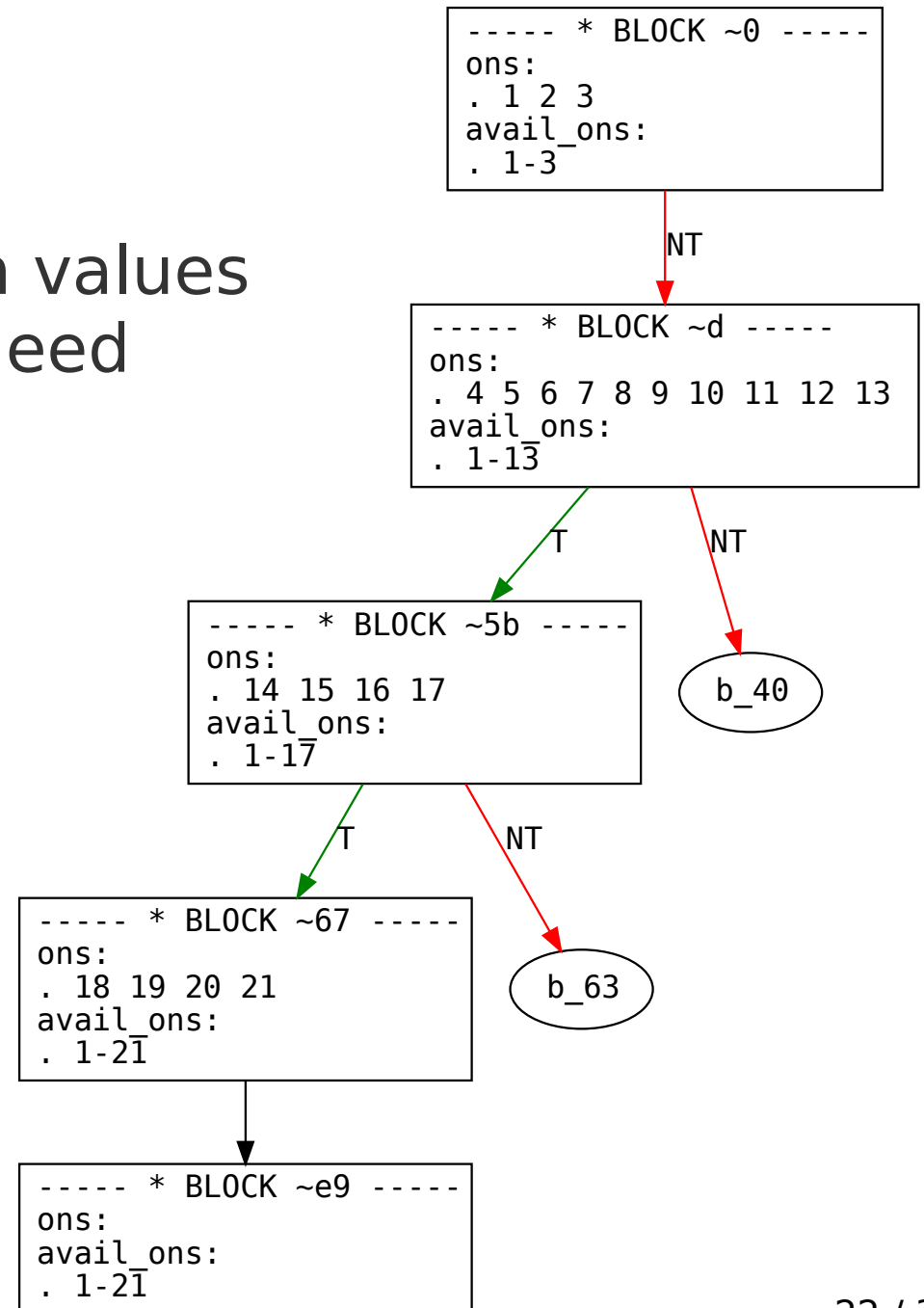
# Analysis  9/12

- In each block, finds which values are available and which need to be computed (common subexpression elimination)

```
----- * BLOCK ~0 -----
ons:
. 1 2 3
avail_ons:
. 1-3
```

NT

```
----- * BLOCK ~d -----
ons:
. 4 5 6 7 8 9 10 11 12 13
avail_ons:
. 1-13
```

T          NT

```
----- * BLOCK ~5b -----
ons:
. 14 15 16 17
avail_ons:
. 1-17
```

b_40

T          NT

```
----- * BLOCK ~67 -----
ons:
. 18 19 20 21
avail_ons:
. 1-21
```

b_63

```
----- * BLOCK ~e9 -----
ons:
avail_ons:
. 1-21
```

# Analysis 10/12

- Calculates the new expressions

```
----- ON CALCS -----
0 = ON_0_RESERVED
1 = #40
2 = #80
3 = MSTORE 0 1 2
4 = #10000...
5 = #0
6 = #6ae17ab7
7 = #5b
8 = #ffffffff
9 = CALLDATALOAD 5
10 = DIV 9 4
11 = AND 10 8
12 = EQ 11 6
13 = JUMPI 7 12
14 = #67
15 = CALLVALUE
16 = ISZERO 15
17 = JUMPI 14 16
18 = #24
```

# Analysis 11/12

- Synthesizes predictor code

```
~0 | ENTRY
    1 = #40
    2 = #80
    3 = MSTORE 0 1 2
~d | ~0
    4 = #10000...
    5 = #0
    6 = #6ae17ab7
    7 = #5b
    8 = #ffffffff
    9 = CALLDATALOAD 5
   10 = DIV 9 4
   11 = AND 10 8
   12 = EQ 11 6
   13 = JUMPI 7 12
....
```

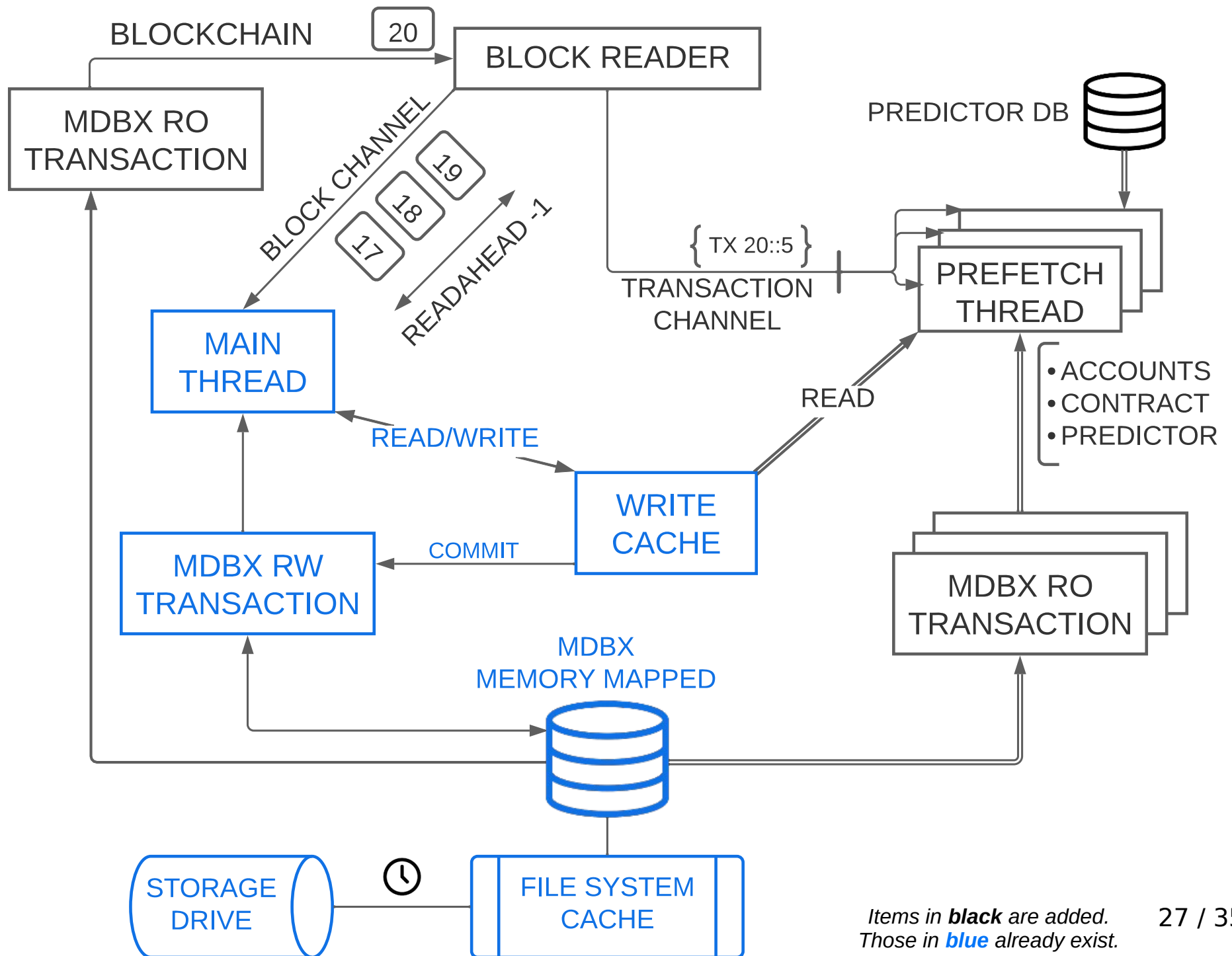# Analysis 12/12

- Outputs in binary encoding, imported in the DB

```
Key = 00a5e63813215d7783df9673e42ec7e1d2e5c0896f17e
96ef6f8d28f1e19f663 (original contract's code hash)

Value = 6a000d001000000100005b00680000010d0067007c0
000015b0079009800000010701e900cb0000016700f400d60000
01cd01fc00d9000001f4000501e4000001cd010701ec0000010
5013401f6000002e900fc00a301980100013401b701c3010001
a301cd01e4010001b701010100400102008090010002002200d0
01d040001000000000000000000000000000000000000000000
000000000000001050000040600...
```

# Prefetching

- Added to erigon, as a separate go package

BLOCKCHAIN  20

BLOCK READER

MDBX RO TRANSACTION

PREDICTOR DB

BLOCK CHANNEL  17 18 19

READAHEAD -1

MAIN THREAD

{ TX 20::5 }

TRANSACTION CHANNEL

PREFETCH THREAD

READ/WRITE

READ

• ACCOUNTS
• CONTRACT
• PREDICTOR

WRITE CACHE

COMMIT

MDBX RW TRANSACTION

MDBX RO TRANSACTION

MDBX MEMORY MAPPED

STORAGE DRIVE

FILE SYSTEM CACHE

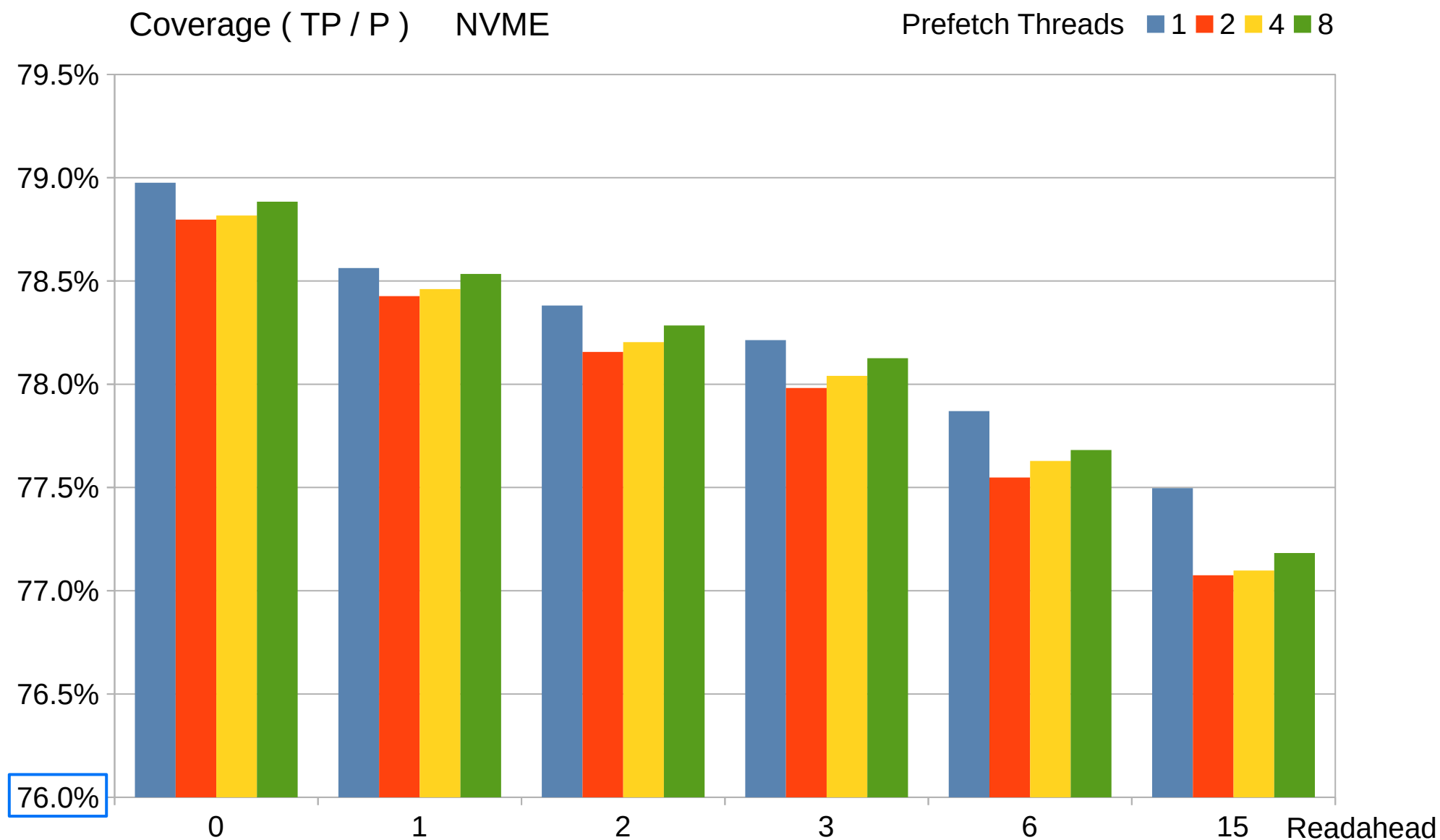*Items in **black** are added.*
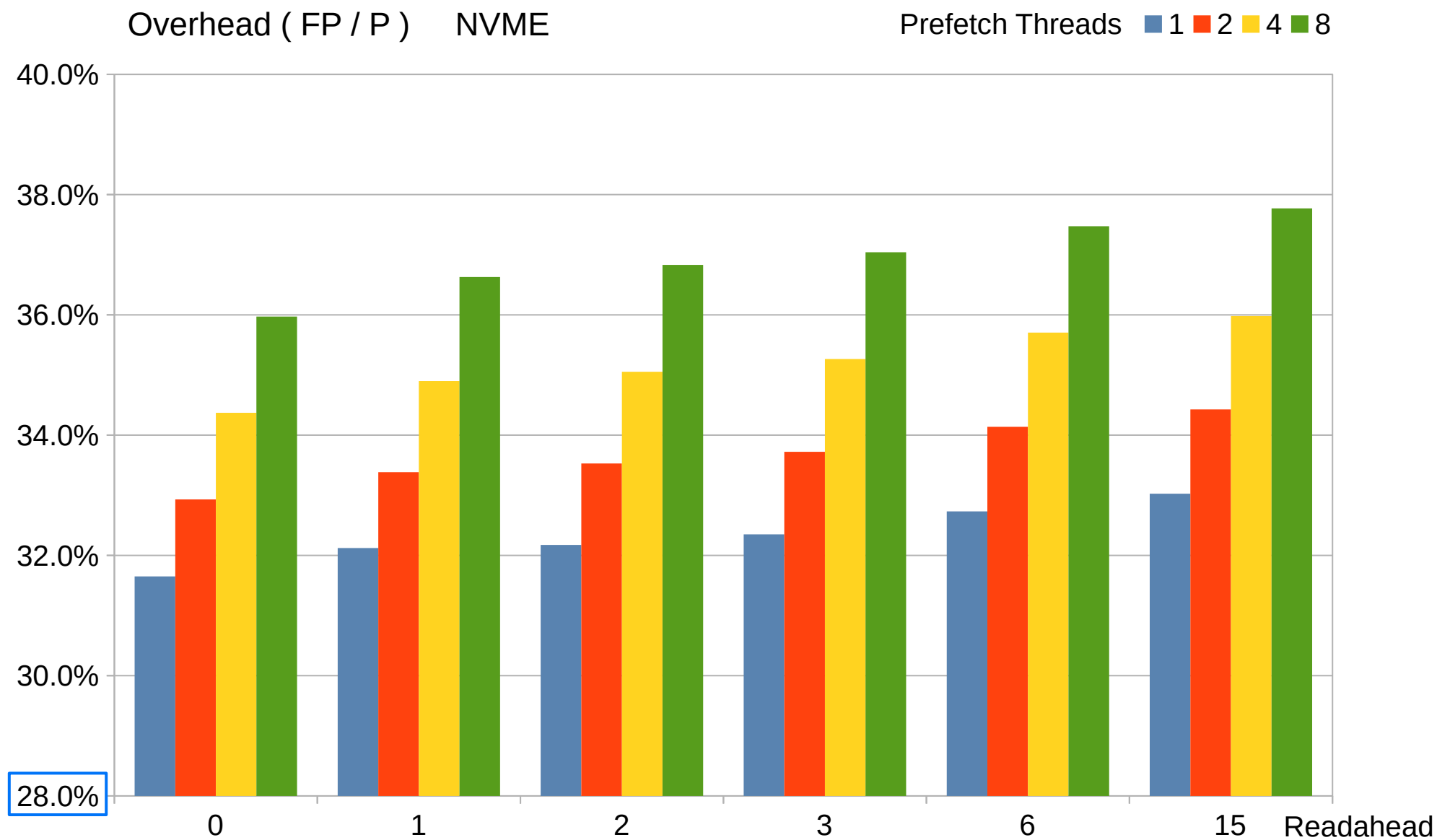*Those in **blue** already exist.*

# Experiments

- All on the same machine

- CPU: zen 2, 3.8 GHz, limited to 1 CCX 3C/6T, shared L3

- RAM: limited by reserving using separate process
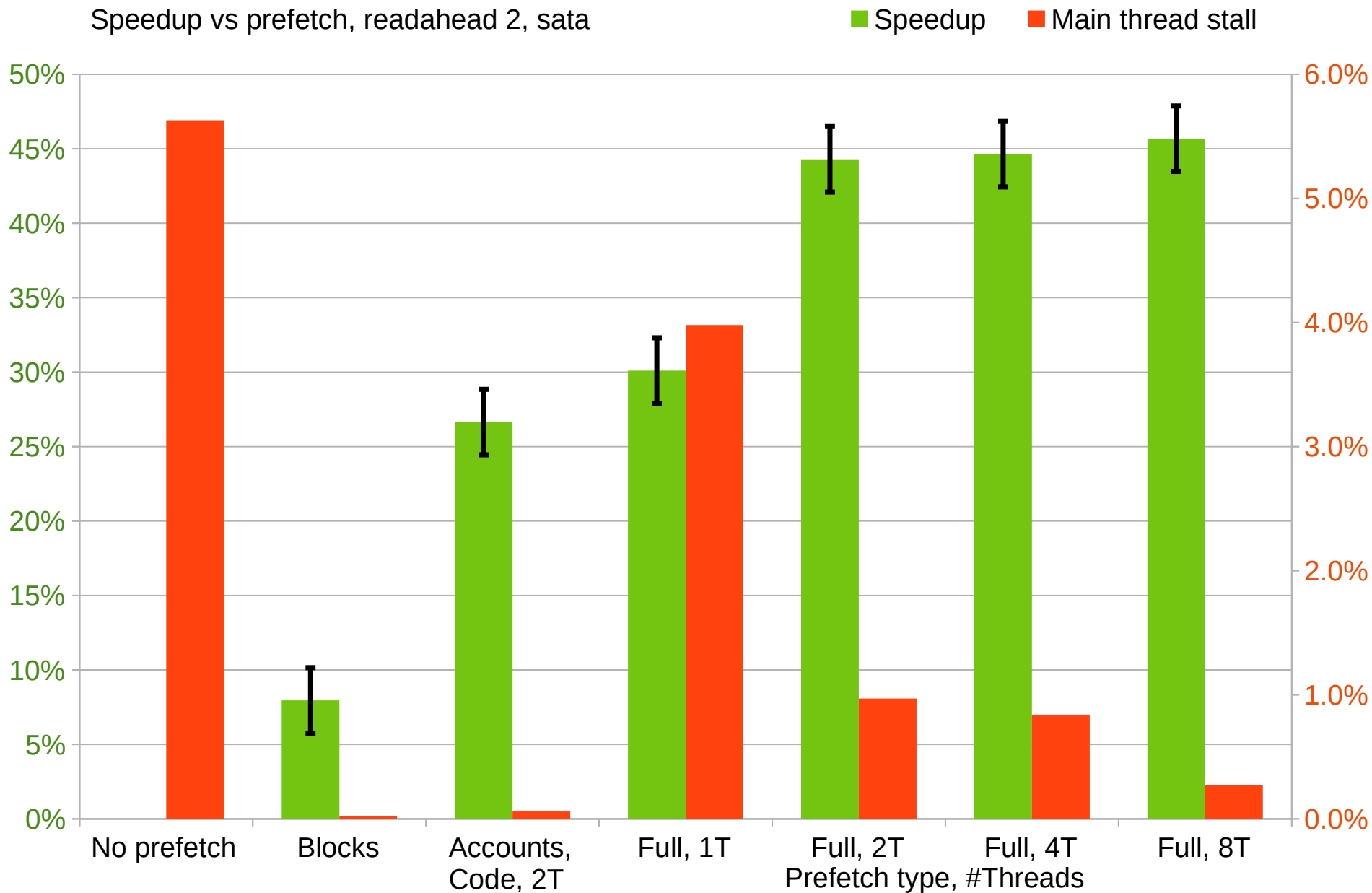
- Disk 1: nvme tlc ssd

- Disk 2: sata qlc ssd

Inter-core one-way data latency between CPU cores

Success % vs readahead

# Overhead % vs readahead



Overhead ( FP / P )    NVME    Prefetch Threads  ■ 1 ■ 2 ■ 4 ■ 8

# Speedup, sata

Speedup vs prefetch, readahead 2, sata

■ Speedup   ■ Main thread stall



Prefetch type, #Threads

| No prefetch | Blocks | Accounts, Code, 2T | Full, 1T | Full, 2T | Full, 4T | Full, 8T |

# Speedup, nvme

Speedup vs prefetch, readahead 2, nvme

■ Speedup    ■ Main thread stall



Prefetch type, #Threads

X-axis categories: No prefetch, Blocks, Accounts, Code, 2T, Full, 1T, Full, 2T, Full, 4T

Left axis: 0%, 5%, 10%, 15%, 20%, 25%

Right axis: 0.0%, 0.5%, 1.0%, 1.5%, 2.0%, 2.5%, 3.0%, 3.5%, 4.0%

# Scatter with all test runs

Blocks / second vs % of time main thread waiting DB



$y = -58.28\,x + 55.91$
$R^2 = 0.95$

NVME, 16GiB, Full prefetch 4T

SATA, 16GiB, No prefetch

SATA, 16GiB, Full prefetch 4T

# Speedup vs block date



Speedup over different time periods, nvme, 16GiB, 4T

Prefetch type: ■ Accounts & Code ◆ Full

Test run block height (x1000)

April 2019

March 2021

# Conclusion

- Significant speedup in an already optimized client

- Close enough to the estimated maximum speedup for a perfect prefetch prediction

- Since it predicts the transactions' read/write set, it effectively determines their dependencies

  → New possibilities for future extensions, by parallelizing the main execution