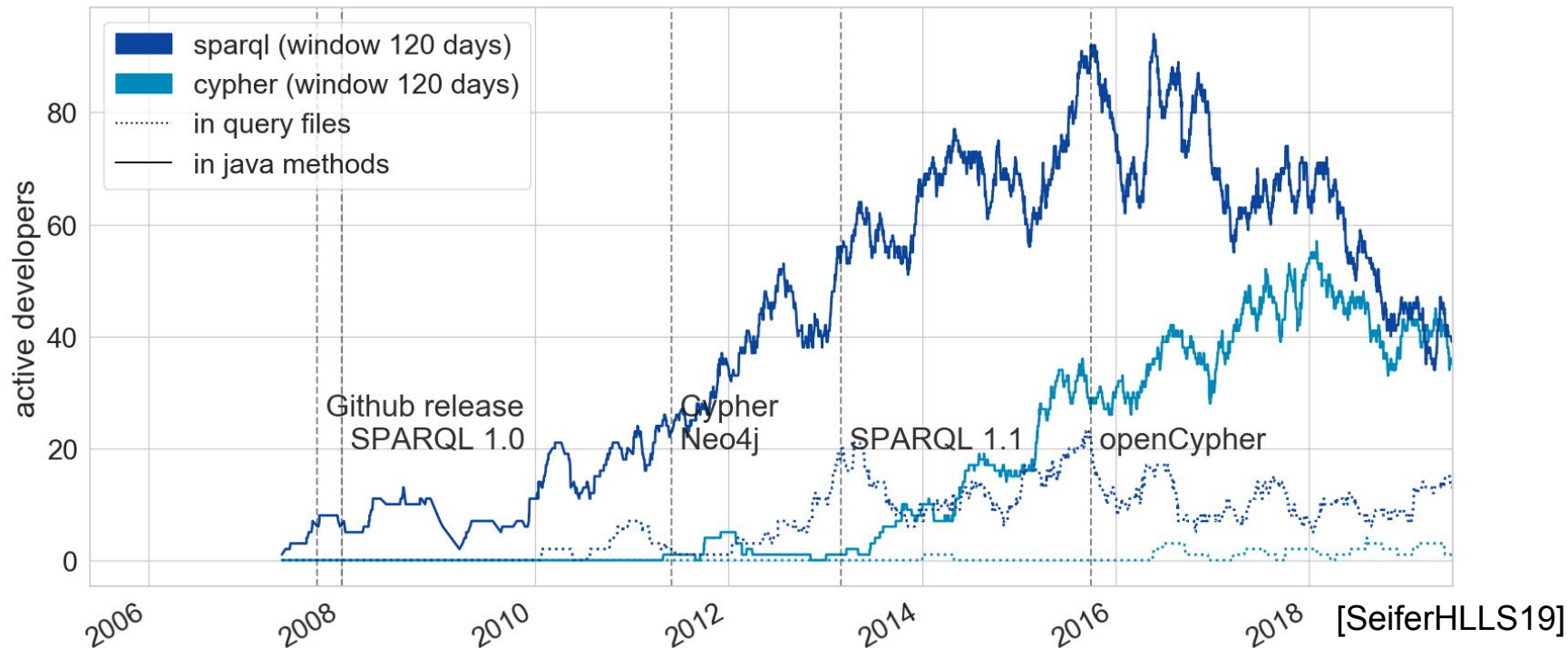


Incremental Map-Reduce on Repository History

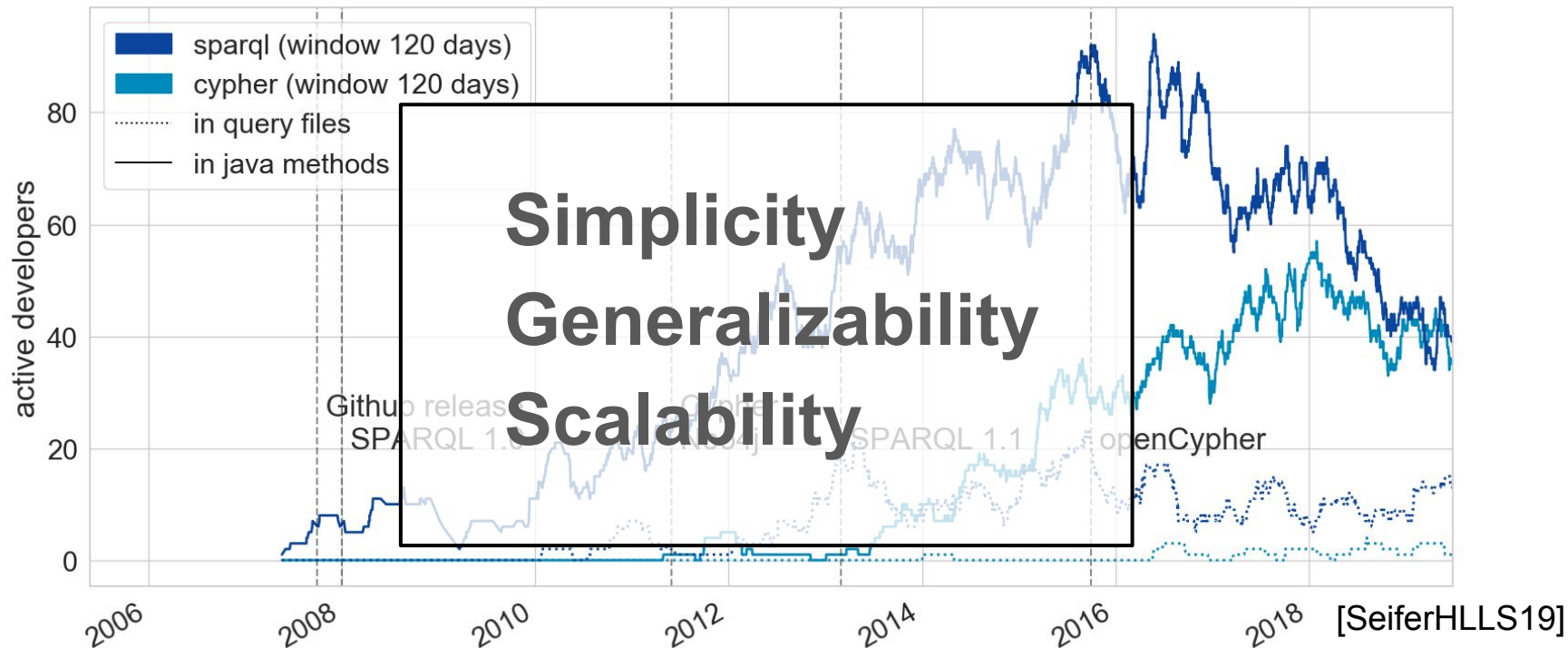
Johannes Härtel and Ralf Lämmel
University of Koblenz-Landau, Germany
Software Languages Team

Challenges of Working with Repository History



7274 repositories, 1.373.000 commits where 22.000 touch sparql/cypher queries

Challenges of Working with Repository History



7274 repositories, 1.373.000 commits where 22.000 touch sparql/cypher queries

Demo

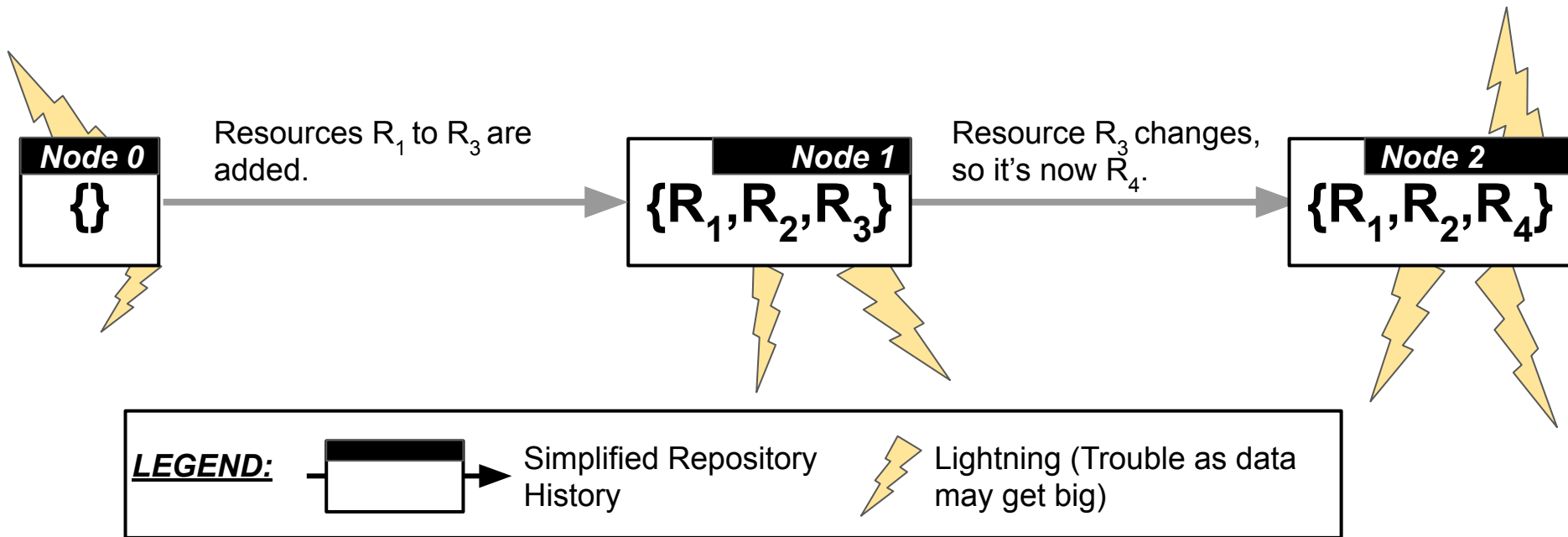
1. How to **represent** History?
2. How to **process** History?
3. How to **memoize** History?

How to represent History?

Use Abelian Groups!

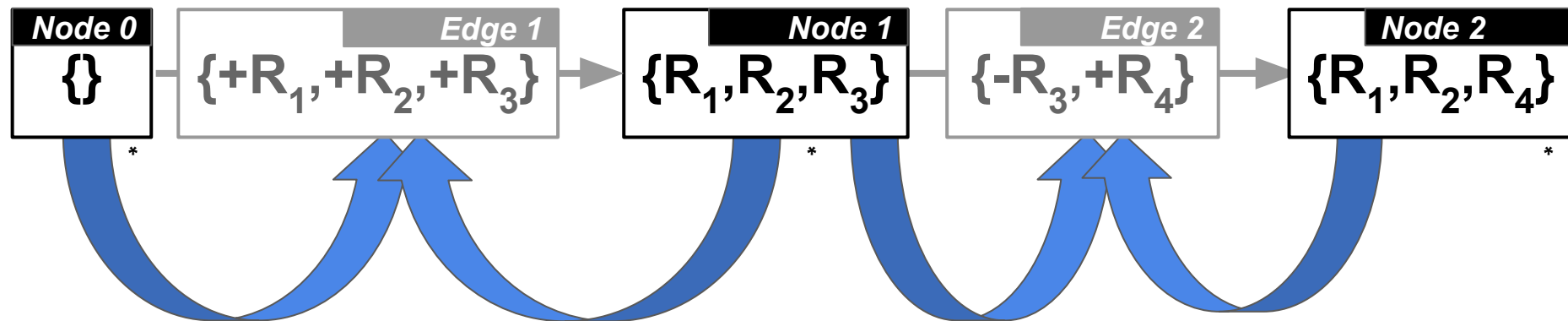
How to represent History?

We model repository history as the Nodes 0-2, storing bags ('{' '}') of immutable resources (R_{1-4}).



How to represent History?

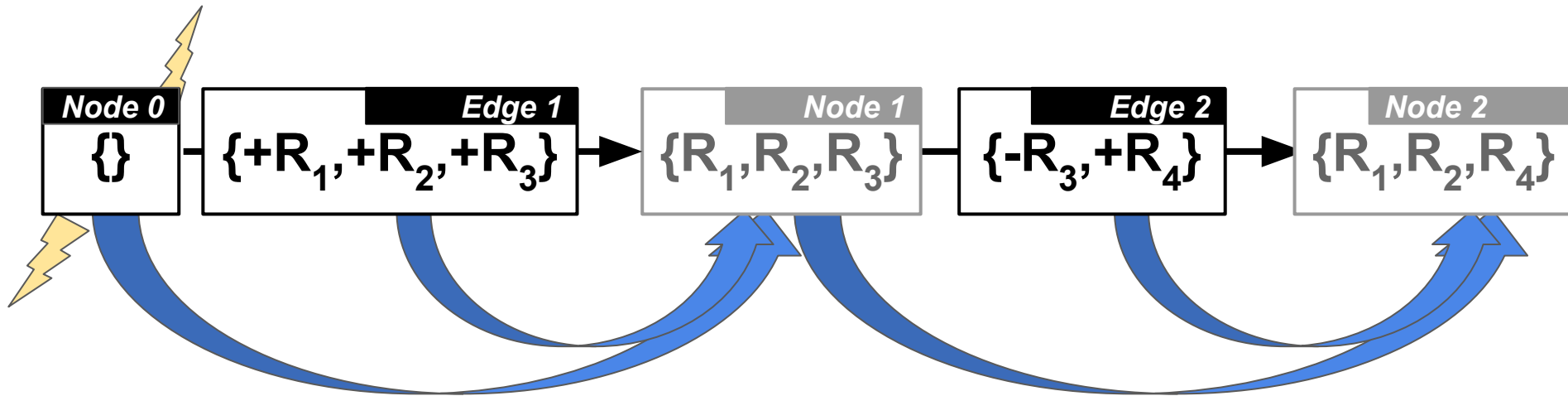
We use Nodes 0-2 to compute the changes on Edges 1-2 in terms of bags with positive and negative elements ($\pm R_{1-4}$).



* *lightnings are avoided due to limitations in space on slide*

How to represent History?

Data on Node 1 and Node 2 can be dropped, as we can use Node 0, Edge 1 and Edge 2 for reproducing all historic data.

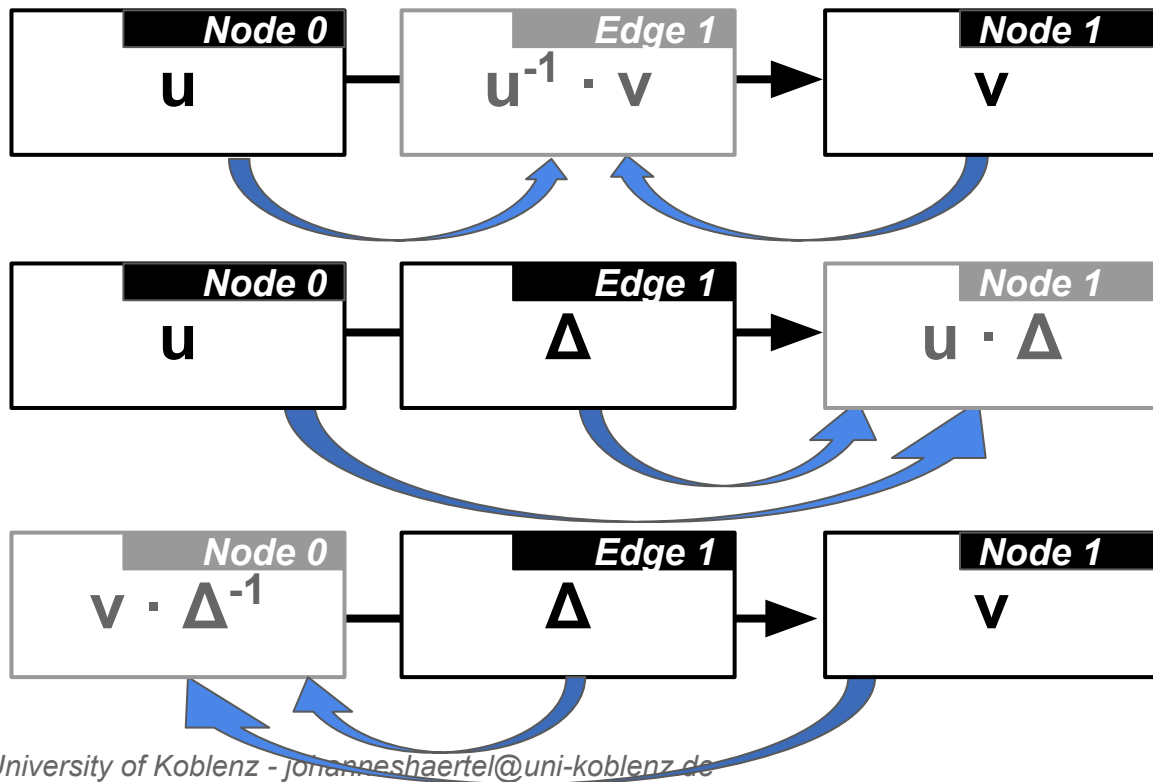


How to represent History?

Use Abelian Groups!

Abelian groups (for u, v and Δ) consist of:

1. A group operator ' \cdot ' (commutative).
2. An inverse element x^{-1} .
3. An identity element z .



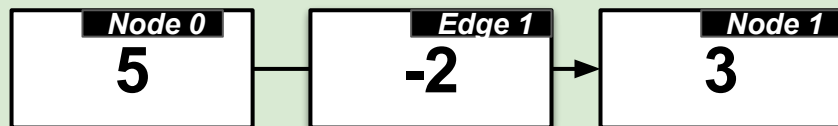
Some Abelian Groups

Integers (Int)

Group operator $+$.

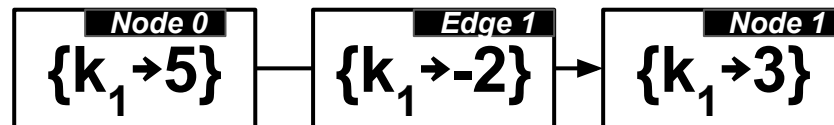
Inverser element $-x$.

Identity element 0 .



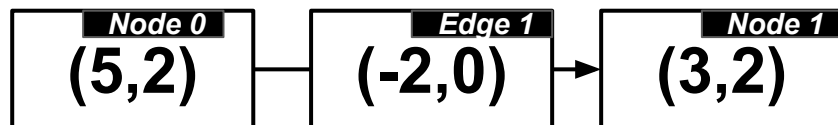
Maps (Map[K, V])

Group operator does an element-wise product, using the operator of a **nested group V**; K can hold arbitrary types . . .



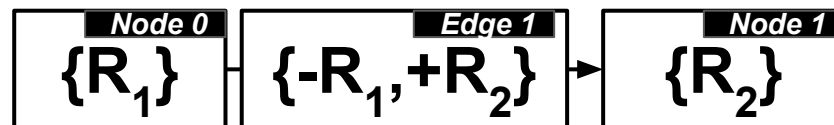
Tuples ((A, B))

All operators delegates to **nested groups** for A and B.



Bags (Bag[K])

Alias for Map[K, Int].



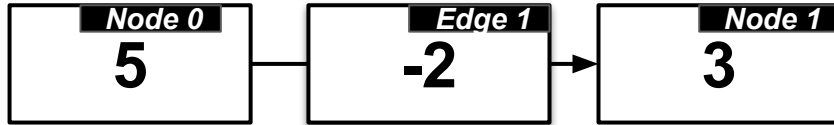
Some Abelian Groups

Integers (Int)

Group operator $+$.

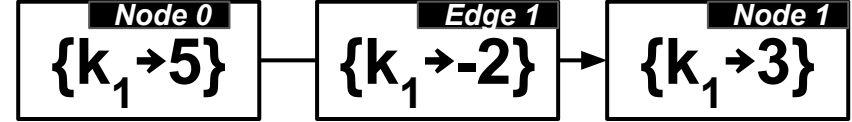
Inverser element $-x$.

Identity element 0 .



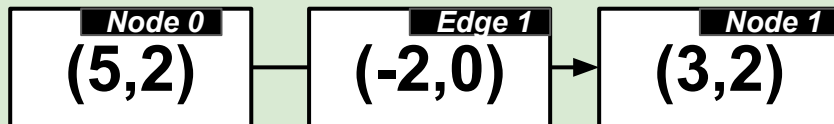
Maps (Map[K, V])

Group operator does an element-wise product, using the operator of a **nested group V**; K can hold arbitrary types . . .



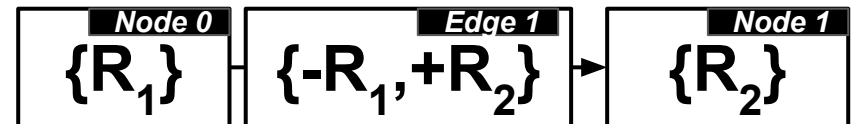
Tuples ((A, B))

All operators delegates to **nested groups** for A and B.



Bags (Bag[K])

Alias for Map[K, Int].



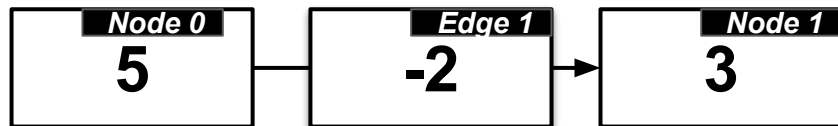
Some Abelian Groups

Integers (Int)

Group operator $+$.

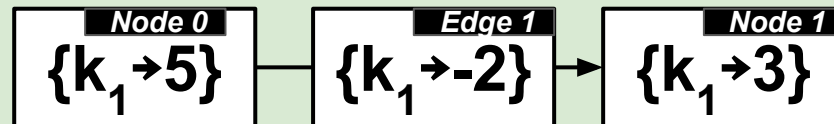
Inverser element $-x$.

Identity element 0 .



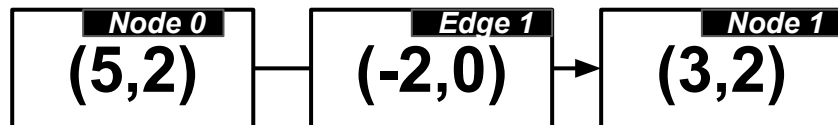
Maps (Map[K, V])

Group operator does an element-wise product, using the operator of a **nested group V**; K can hold arbitrary types . . .



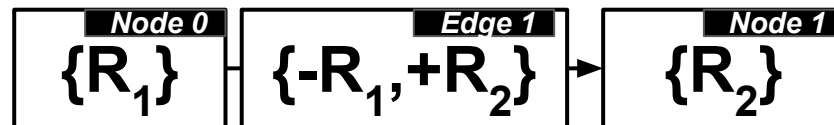
Tuples ((A, B))

All operators delegates to **nested groups** for A and B.



Bags (Bag[K])

Alias for Map[K, Int].



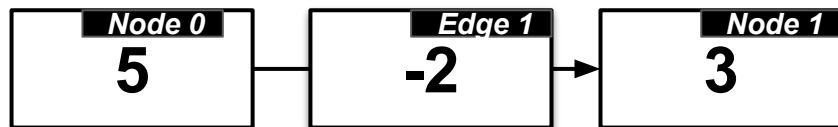
Some Abelian Groups

Integers (Int)

Group operator $+$.

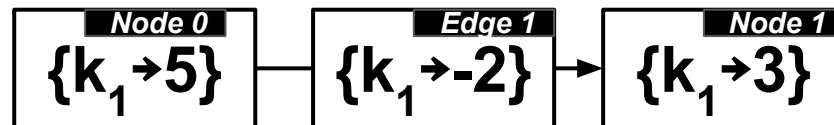
Inverser element $-x$.

Identity element 0 .



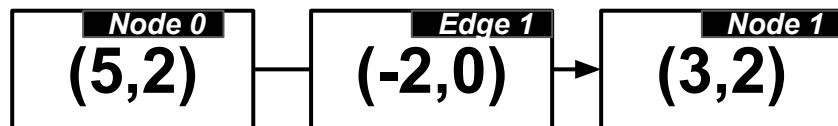
Maps (Map[K, V])

Group operator does an element-wise product, using the operator of a **nested group V**; K can hold arbitrary types . . .



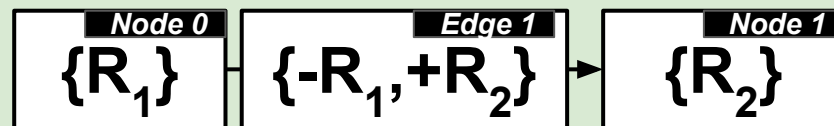
Tuples ((A, B))

All operators delegates to **nested groups** for A and B.



Bags (Bag[K])

Alias for Map[K, Int].

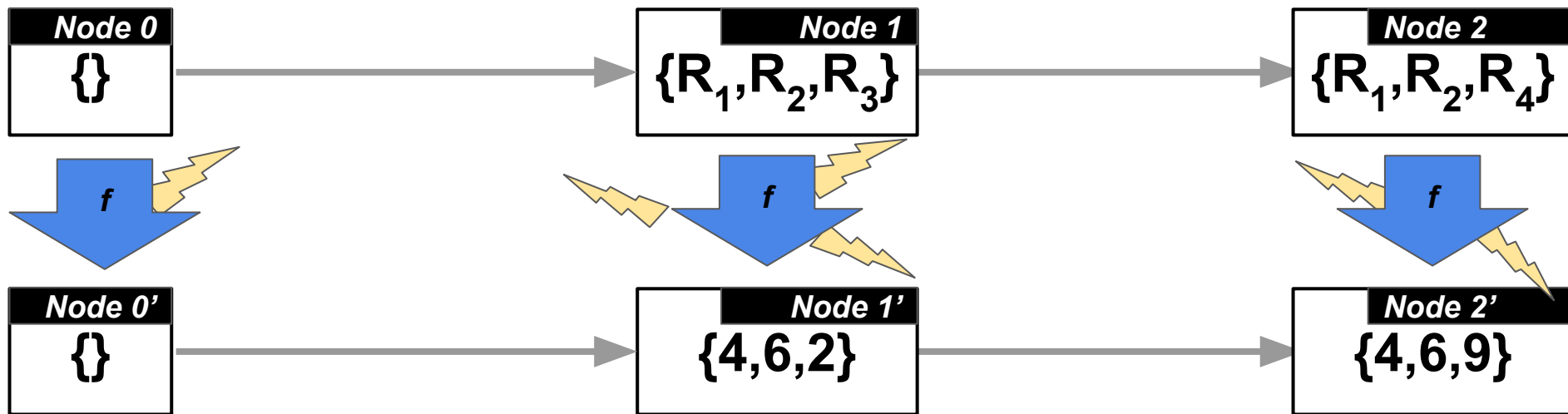


How to process History?

Use Group Homomorphisms!

How to process History?

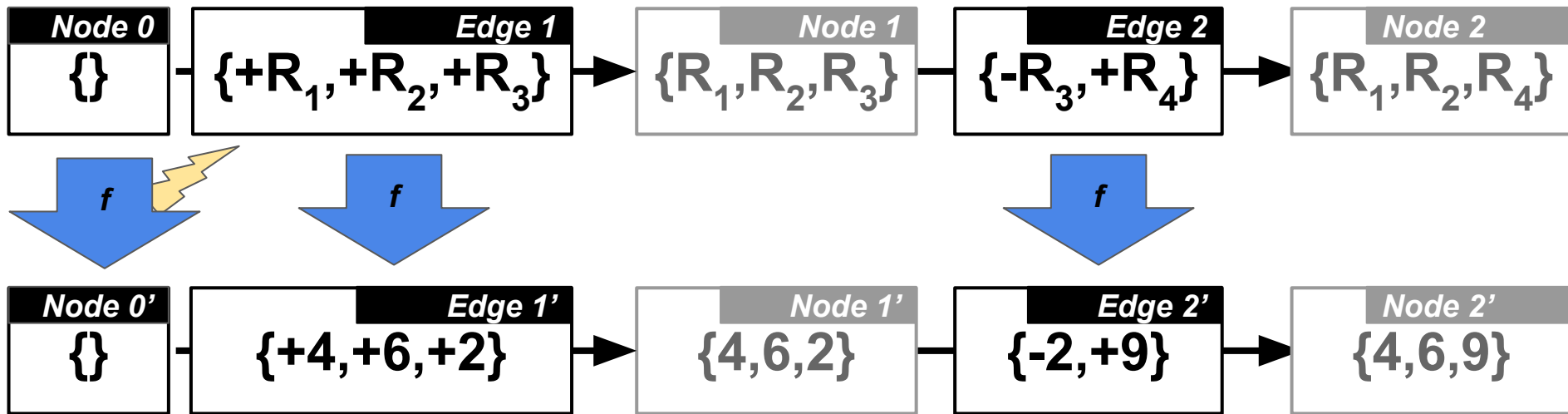
Extracting metrics can be done by applying a function f to the bags on Node 0-2.



(In this example metrics are $f(\{R_1\}) = 4$, $f(\{R_2\}) = 6$, $f(\{R_3\}) = 2$, $f(\{R_4\}) = 9$)

How to process History?

We can apply f to Node 0, Edge 1 and Edge 2 and get the same results;
Node 1' and Node 2' can be restored.

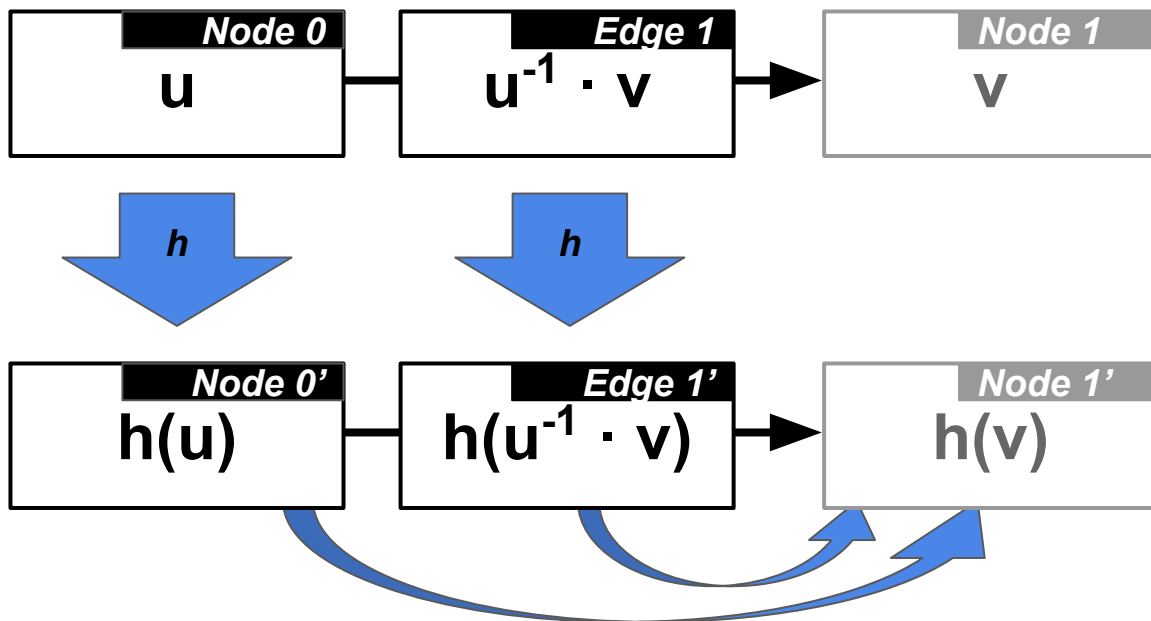


(In this example metrics are $f(\{R_1\}) = 4$, $f(\{R_2\}) = 6$, $f(\{R_3\}) = 2$, $f(\{R_4\}) = 9$)

How to process History?

Use Group Homomorphisms!

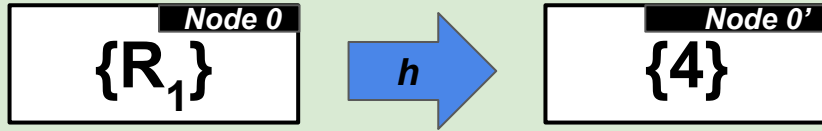
For a group homomorphism h ,
 $h(u^{-1} \cdot v) = h(u)^{-1} \cdot h(v)$ holds.



Some Group Homomorphisms

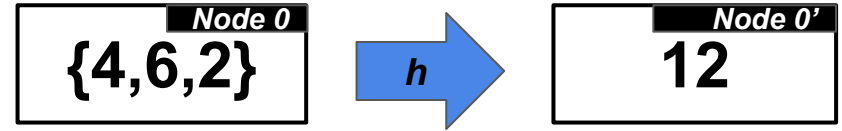
map (Bag[V1] → Bag[V2])

Maps the elements of a bag.



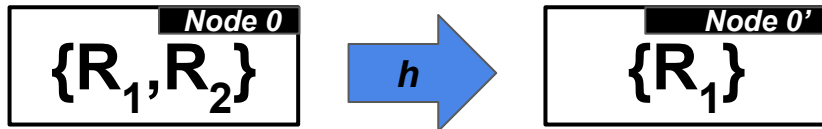
sum (Bag[Int] → Int)

Sums up the integer entries of a bag producing an integer.



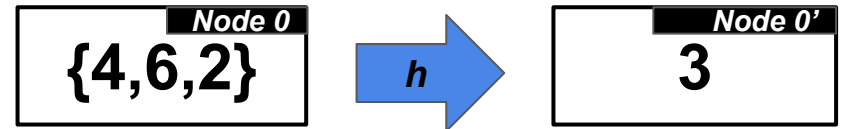
filter (Bag[V1] → Bag[V2])

Filters the elements of a bag.



count (Bag[K] → Int)

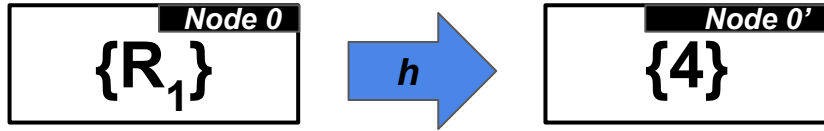
Counts the elements of a bag producing an integer.



Some Group Homomorphisms

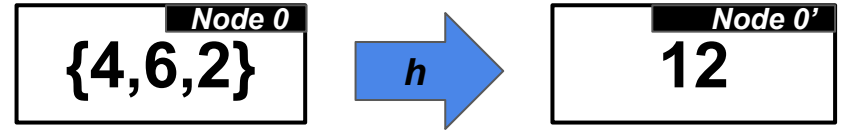
map (Bag[V1] → Bag[V2])

Maps the elements of a bag.



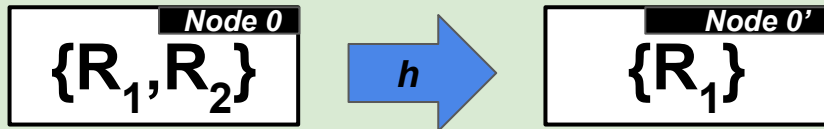
sum (Bag[Int] → Int)

Sums up the integer entries of a bag producing an integer.



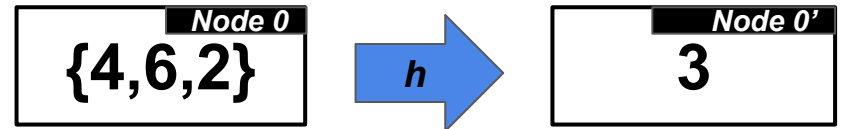
filter (Bag[V1] → Bag[V2])

Filters the elements of a bag.



count (Bag[K] → Int)

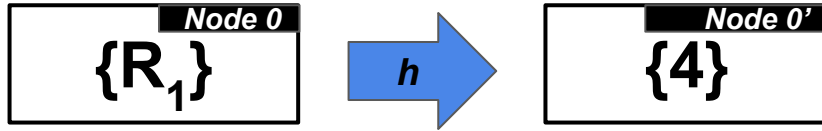
Counts the elements of a bag producing an integer.



Some Group Homomorphisms

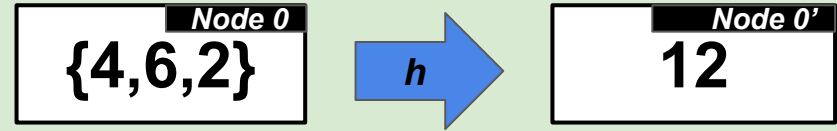
map (Bag[V1] → Bag[V2])

Maps the elements of a bag.



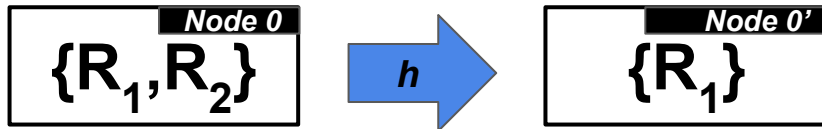
sum (Bag[Int] → Int)

Sums up the integer entries of a bag producing an integer.



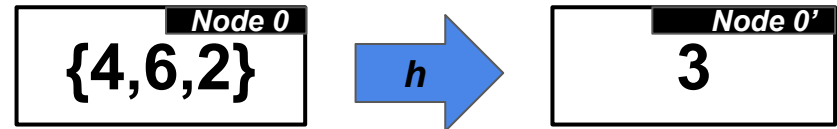
filter (Bag[V1] → Bag[V2])

Filters the elements of a bag.



count (Bag[K] → Int)

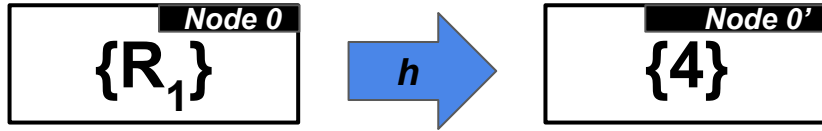
Counts the elements of a bag producing an integer.



Some Group Homomorphisms

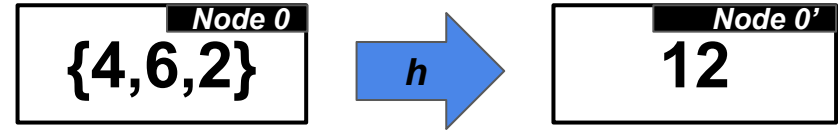
map (Bag[V1] → Bag[V2])

Maps the elements of a bag.



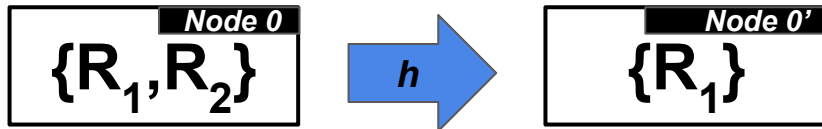
sum (Bag[Int] → Int)

Sums up the integer entries of a bag producing an integer.



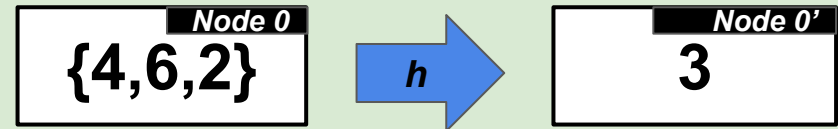
filter (Bag[V1] → Bag[V2])

Filters the elements of a bag.



count (Bag[K] → Int)

Counts the elements of a bag producing an integer.

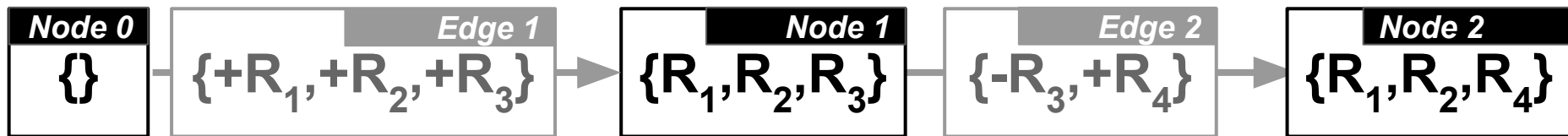


How to memoize History?

Invert it!

How to memoize History?

Instead of pointing from nodes/edges to the data, history can be inverted, pointing from data to the edges/nodes.

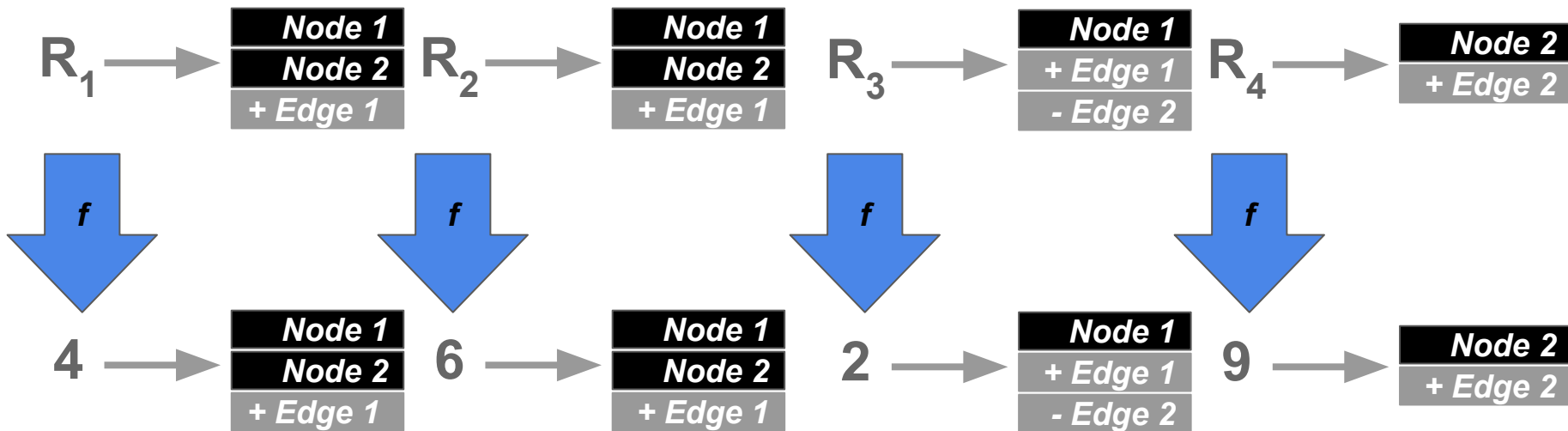


INVERT



How to memoize History?

Functions can be applied without causing redundant computations.



(Metrics are $f(\{R_1\}) = 4$, $f(\{R_2\}) = 6$, $f(\{R_3\}) = 2$, $f(\{R_4\}) = 9$)

Evaluation

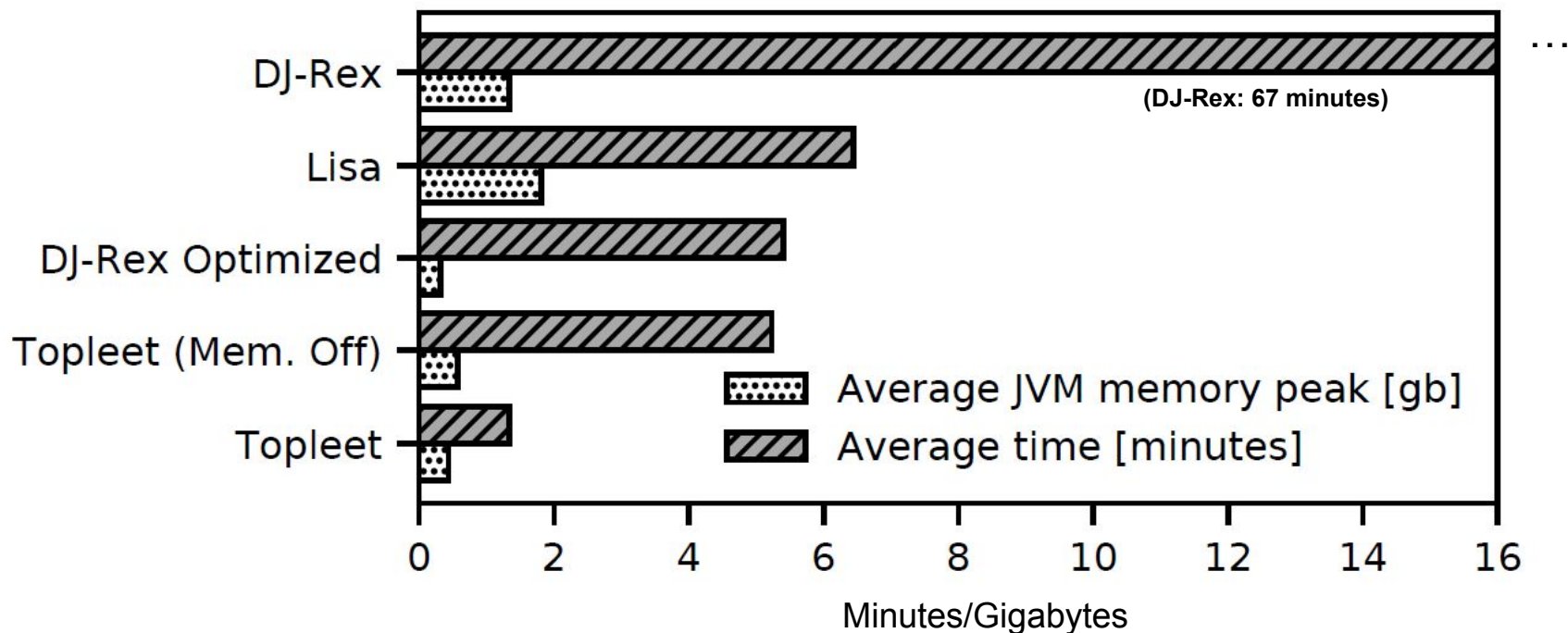
* Excluded from evaluation.

Evaluated Approaches

	Language/ Interface	History Model	Distribution	Parallelization	Compression/ Red. Reduction/ Memoization	Change Processing/ Incrementalization
DJ-Rex [ShangJAH09]	Map-Reduce		Generic	Generic	Manual Optimization	Manual Optimization
Boa* [DyerNRN13]	DSL		Generic	Generic	Specific to Java	Manual Optimization
LISA [AlexandruG15]	Signal/Collect	Linear		Specific to Resources	Specific to AST	Incremental Messaging
Topleet	Map-Reduce	Acyclic	Generic	Generic	Abelian groups	Group Homomorphisms & Inc. Indexing

Average Time and Memory Usage

Running cyclomatic complexity solutions on 98 repositories.



Try Topleet

<https://github.com/topleet/topleet>

- **Extended map-reduce functionality** is provided, such as, join or cartesian.
- **Concise interface** defined by 5 core operators.
- **Interchangeable implementations**, e.g., incremental vs. non-incremental.
- **Differential testing** is used for checking correctness.
- **More application** cases are available online.
- **Future work** may include authorship attribution, refactoring detection and graph queries to the history.

Fin

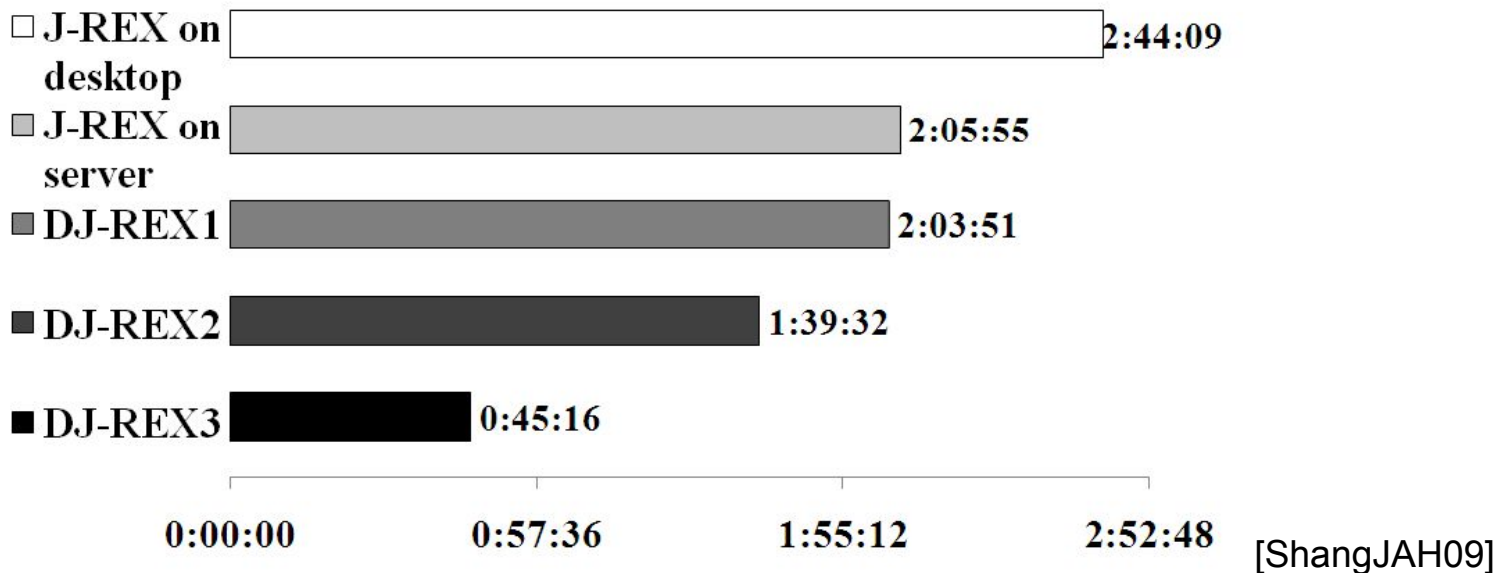
<https://github.com/topleet/topleet>

References

- [ShangJAH09] Shang, Weiyi, Zhen Ming Jiang, Bram Adams, and Ahmed E. Hassan. "MapReduce as a general framework to support research in Mining Software Repositories (MSR)." In 2009 6th IEEE International Working Conference on Mining Software Repositories, pp. 21-30. IEEE, 2009.
- [DyerNRN13] Dyer, R., Nguyen, H. A., Rajan, H., & Nguyen, T. N. (2013, May). Boa: A language and infrastructure for analyzing ultra-large-scale software repositories. In 2013 35th International Conference on Software Engineering (ICSE) (pp. 422-431). IEEE.
- [StutzBC10] Stutz, Philip, Abraham Bernstein, and William Cohen. "Signal/collect: graph algorithms for the (semantic) web." In International Semantic Web Conference, pp. 764-780. Springer, Berlin, Heidelberg, 2010.
- [AlexandruG15] Alexandru, Carol V., and Harald C. Gall. "Rapid multi-purpose, multi-commit code analysis." In 2015 IEEE/ACM 37th IEEE International Conference on Software Engineering, vol. 2, pp. 635-638. IEEE, 2015.
- [SeiferHLLS19] Seifer, Philipp, Härtel, Johannes, Leinberger, Martin, Lämmel, Ralf and Staab, Steffen (2019) Empirical study on the usage of graph query languages in open source Java projects. 12th ACM SIGPLAN International Conference on Software Language Engineering (SLE 2019), Greece. 20 - 25 Oct 2019. 15 pp .

DJ-Rex Solution: Migration to Distributed Map-Reduce

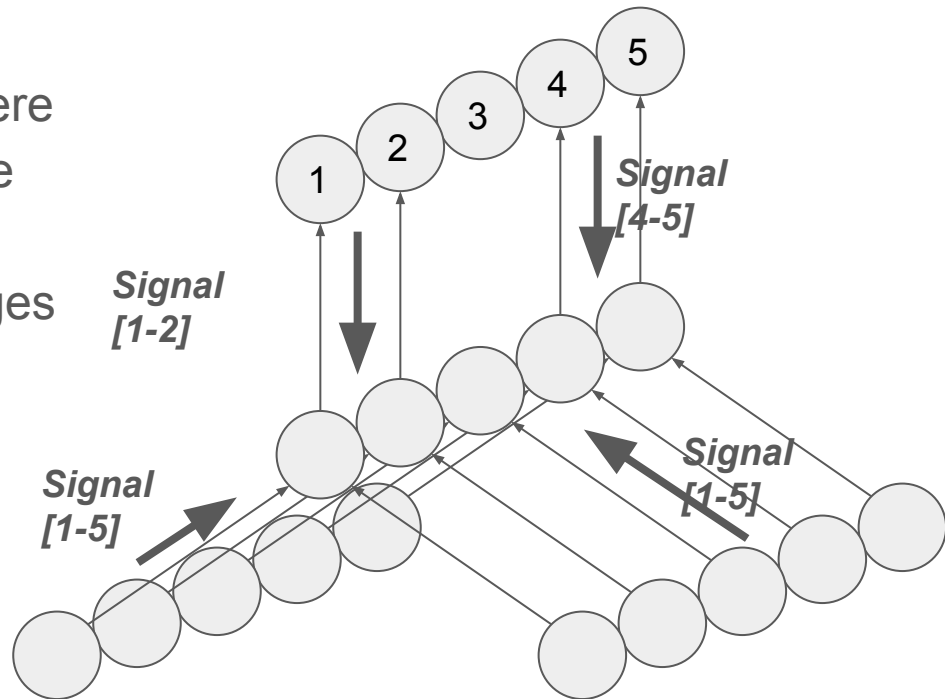
Successively adding Hardware leads to runtime improvements (DJ-REX1-3).



LISA Solution:

Reduction of Redundancies

LISA uses Signal/Collect [StutzBC10] to express queries, i.e., an actor system where asynchronous messages are sent until the system converges. LISA creates a multi-revision AST and just sends messages once for successive but equal version ranges.



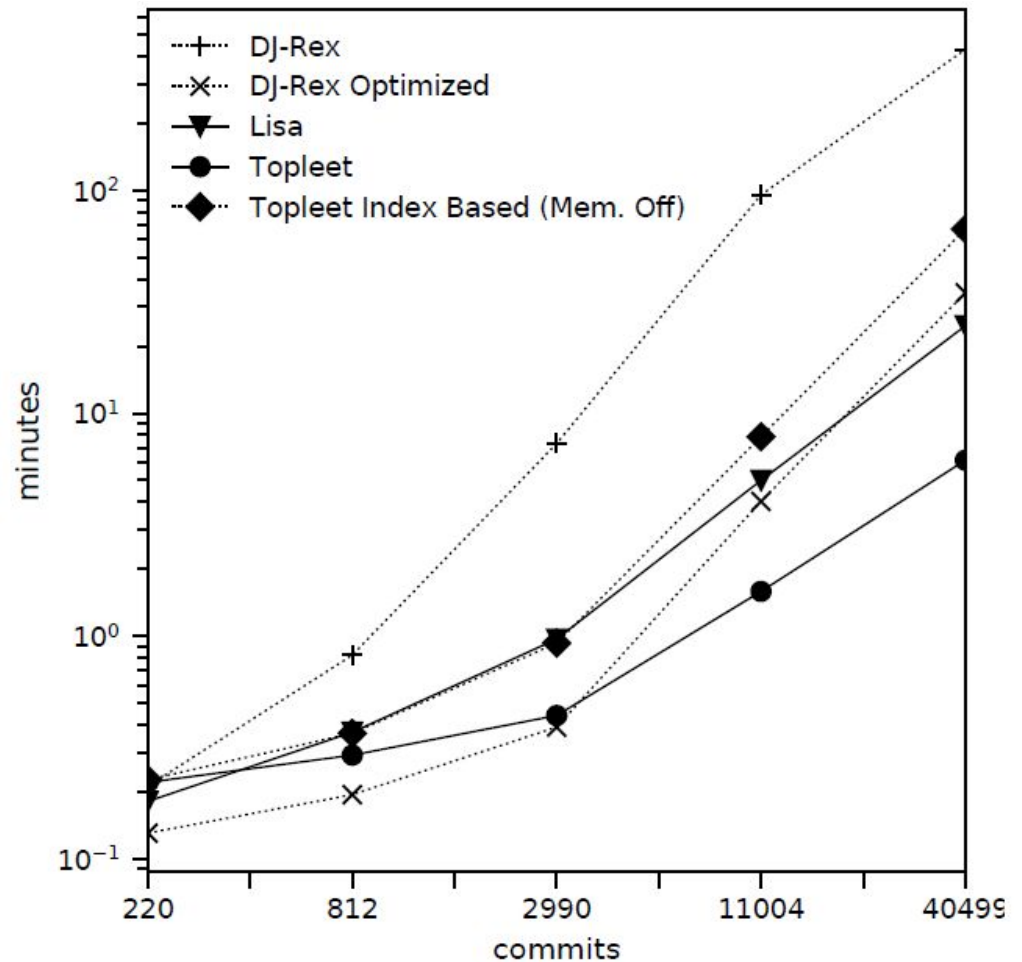
Boa Solution: Domain Specific Languages

Analysis is written in a DSL and compiled to Map-Reduce. The right except answers: “How many fixing revisions added null checks?” (copy form the documentation). It contains handling of changed resources.

```
before node: ChangedFile -> {  
  # if this is a fixing revision and  
  # there was a previous version of the file  
  if (isfixing && haskey(files, node.name)) {  
    # count how many null checks were previously in the file  
    count = 0;  
    visit(getast(files[node.name]));  
    last := count;  
  
    # count how many null checks are currently in the file  
    count = 0;  
    visit(getast(node));  
  
    # if there are more null checks, output  
    if (count > last)  
      AddedNullCheck << 1;  
  }  
  if (node.change == ChangeKind.DELETED)  
    remove(files, node.name);  
  else  
    files[node.name] = node;  
  stop;  
}
```

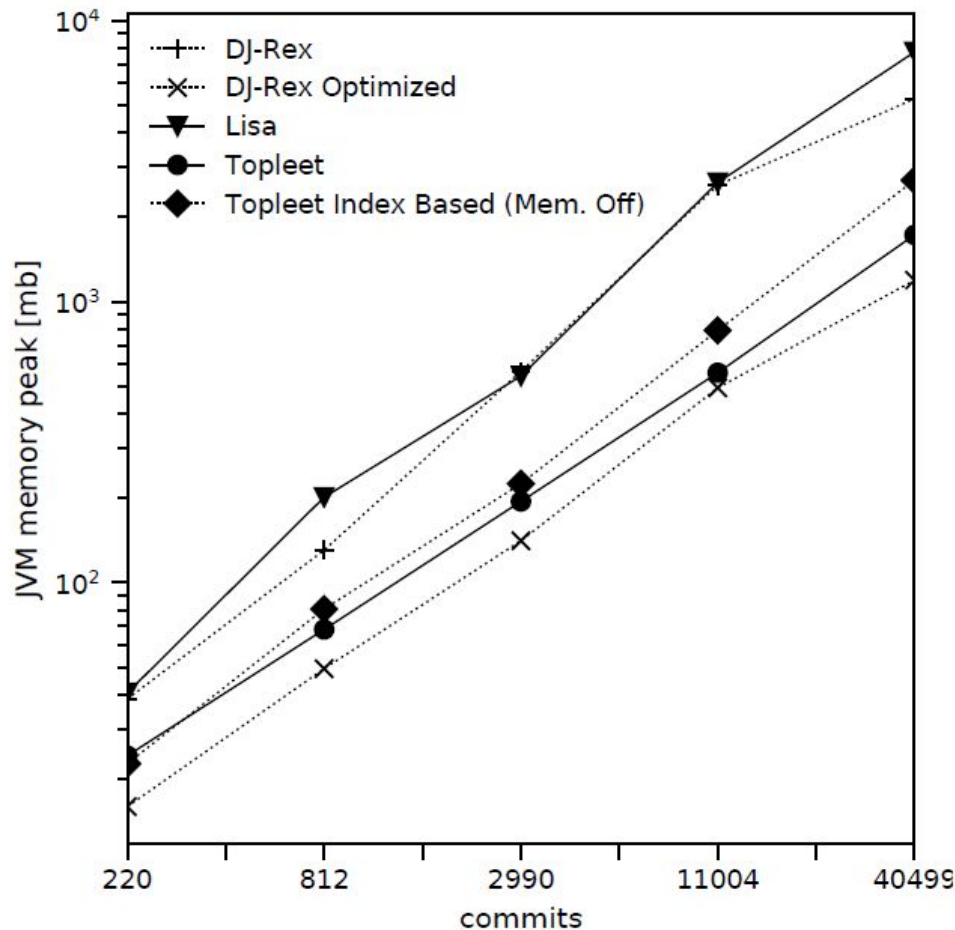
Time vs. Commits

Running the cyclomatic complexity solutions on 98 repositories grouped by commit count.



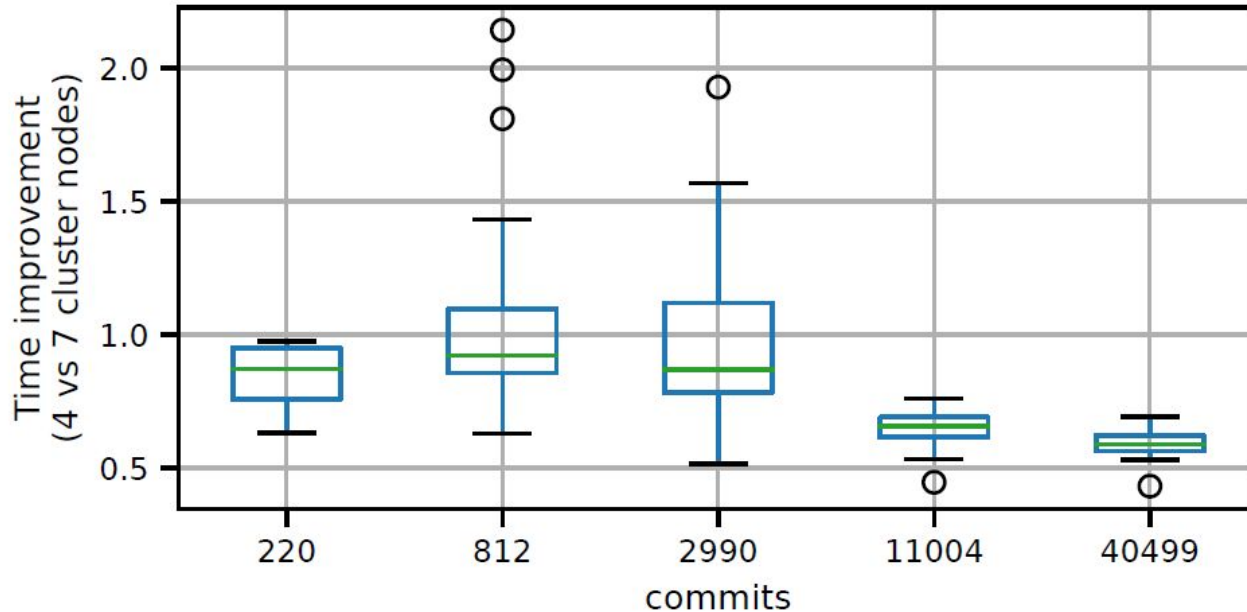
Memory vs. Commits

Running the cyclomatic complexity solutions on 98 repositories grouped by commit count.



Distribution vs. Commits

Running Topleet cyclomatic complexity in distribution on all 98 repositories grouped by commit count.



Demo Backup

Demo

```
9  ▶  object Sandbox {  
10  
11    // The foreign function computing MCC on a resource.  
12    def computeMCC(resource: Resource): Int = {...}  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71    // Library imports.  
72    import org.topleet.libs.Natives._  
73    import org.topleet.libs.Gits._  
74  
75    implicit val engine: Engine = IncrementalParallelEngine.create()  
76  
77    ▶  def main(args: Array[String]): Unit = {  
78      |  
79    }  
80  }  
81
```

Demo

```
9 ▶ object Sandbox {
```

```
10
```

```
11 // The foreign function computing MCC on a resource.
```

```
12 def computeMCC(resource: Resource): Int = {...}
```

```
13
```

```
14 // Library imports.
```

```
15 import org.toplect.libc.Native
```

```
16 import org.toplect.libc.Native
```

```
17
```

```
18 implicit val env = Env()
```

```
19
```

```
20 def main(args: List[String]): Int = {
```

```
21
```

```
22 |
```

```
23
```

```
24 }
```

```
25 }
```

```
26 }
```

```
val parser = ASTParser.newParser(AST.JLS13)
parser.setSource(resource.read().toCharArray)
parser.setKind(ASTParser.K_COMPILATION_UNIT)
```

```
val cu = parser.createAST(monitor = null).asInstanceOf[CompilationUnit]
```

```
var result = 0
```

```
cu.accept(new ASTVisitor() {
```

```
  override def visit(node: WhileStatement): Boolean = {
    result = result + 1
    super.visit(node)
  }
})
```

```
...
```


Demo

```
9  ▶  object Sandbox {  
10  
11    // The foreign function computing MCC on a resource.  
12    def computeMCC(resource: Resource): Int = {...}  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71    // Library imports.  
72    import org.topleet.libs.Natives._  
73    import org.topleet.libs.Gits._  
74  
75    implicit val engine: Engine = IncrementalParallelEngine.create()  
76  
77    def main(args: Array[String]): Unit = {  
78      val shas = git( address = "jwtk/jjwt")  
79    }  
80  }  
81
```

Demo

```
9  ▶  object Sandbox {  
10  
11    // The foreign function computing MCC on a resource.  
12    def computeMCC(resource: Resource): Int = {...}  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71    // Library imports.  
72    import org.topleet.libs.Natives._  
73    import org.topleet.libs.Gits._  
74  
75    implicit val engine: Engine = IncrementalParallelEngine.create()  
76  
77    def main(args: Array[String]): Unit = {  
78      val shas = git( address = "jwtk/jjwt")  
79  
80      shas.show()  
81    }  
82  }  
83
```

Demo

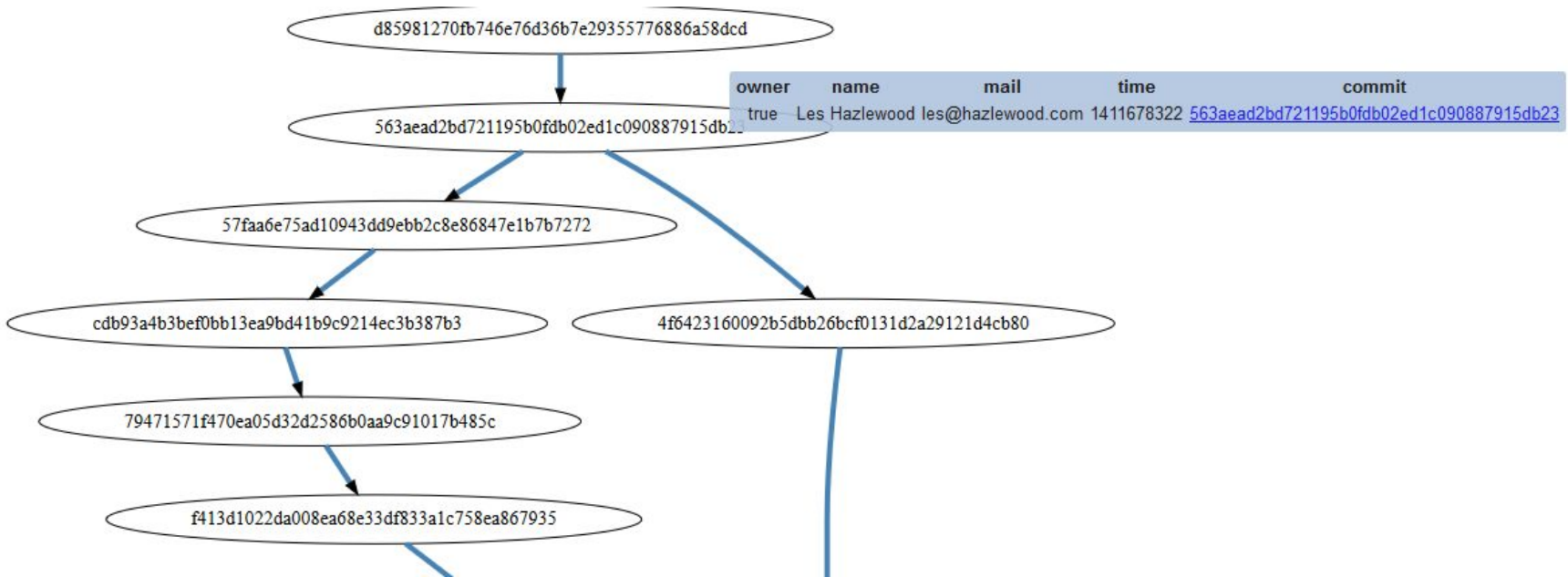
qegal - Google Drive

SANER 20 (+) - Google Slides

Telegram Web

/C:/Data/Repos/topleet/temp/view

file:///C:/Data/Repos/topleet/temp/viewer/viewer.html



Demo

```
9  ▶  object Sandbox {  
10  
11    // The foreign function computing MCC on a resource.  
12    def computeMCC(resource: Resource): Int = {...}  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71    // Library imports.  
72    import org.toplevel.libs.Natives._  
73    import org.toplevel.libs.Gits._  
74  
75    implicit val engine: Engine = IncrementalParallelEngine.create()  
76  
77    def main(args: Array[String]): Unit = {  
78      val shas = git( address = "jwtk/jjwt")  
79    }  
80  }  
81
```

Demo

```
9  ▶  object Sandbox {
10
11      // The foreign function computing MCC on a resource.
12      def computeMCC(resource: Resource): Int = {...}
13
14      // Library imports.
15      import org.topleet.libs.Natives._
16      import org.topleet.libs.Gits._
17
18      implicit val engine: Engine = IncrementalParallelEngine.create()
19
20      def main(args: Array[String]): Unit = {
21          val shas: Leet[SHA, Single[SHA]] = git( address = "jwtk/jjwt")
22      }
23  }
```

Demo

```
9  ▶  object Sandbox {
10
11      // The foreign function computing MCC on a resource.
12      def computeMCC(resource: Resource): Int = {...}
13
14      70
15
16      71      // Library imports.
17
18      72      import org.topleet.libs.Natives._
19      73      import org.topleet.libs.Gits._
20
21      74
22      75      implicit val engine: Engine = IncrementalParallelEngine.create()
23
24      76
25      77      ▶  def main(args: Array[String]): Unit = {
26      78          val shas: Leet[SHA, Single[SHA]] = git( address = "jwtk/jjwt")
27      79
28      80          shas.res
29
30      81      resources()      Leet[SHA, Natives.Bag[(Gits.Path, Resource)]]
31      82
32      83
33      84      Press Ctrl+. to choose the selected (or first) suggestion and insert a dot afterwards: Next Tip
```

Demo

```
9  ▶  object Sandbox {  
10  
11    // The foreign function computing MCC on a resource.  
12    def computeMCC(resource: Resource): Int = {...}  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71    // Library imports.  
72    import org.topleet.libs.Natives._  
73    import org.topleet.libs.Gits._  
74  
75    implicit val engine: Engine = IncrementalParallelEngine.create()  
76  
77    def main(args: Array[String]): Unit = {  
78      val shas: Leet[SHA, Single[SHA]] = git( address = "jwtk/jjwt")  
79  
80      val resources: Leet[SHA, Bag[(Path, Resource)]] = shas.resources()  
81  
82    }  
83  }
```


Demo

```
9  ▶  object Sandbox {
10
11      // The foreign function computing MCC on a resource.
12      def computeMCC(resource: Resource): Int = {...}
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71      // Library imports.
72      import org.topleet.libs.Natives._
73      import org.topleet.libs.Gits._
74
75      implicit val engine: Engine = IncrementalParallelEngine.create()
76
77      ▶  def main(args: Array[String]): Unit = {
78          val shas: Leet[SHA, Single[SHA]] = git( address = "jwtk/jjwt")
79
80          val resources: Leet[SHA, Bag[(Path, Resource)]] = shas.resources()
81
82          val mcCabe: Leet[SHA, Bag[Int]] = resources
83              .filter { case (path, _) => path.endsWith(".java") }
84              .map { case (_, resource) => computeMCC(resource) }
85
86      }
87  }
88
```


Demo

```
9  ▶  object Sandbox {
10
11      // The foreign function computing MCC on a resource.
12      def computeMCC(resource: Resource): Int = {...}
13
14
15
16
17
18
19
20
21      // Library imports.
22      import org.topleet.libs.Natives._
23      import org.topleet.libs.Gits._
24
25
26
27
28      implicit val engine: Engine = IncrementalParallelEngine.create()
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77  ▶  def main(args: Array[String]): Unit = {
78      val shas: Leet[SHA, Single[SHA]] = git( address = "jwtk/jjwt")
79
80      val resources: Leet[SHA, Bag[(Path, Resource)]] = shas.resources()
81
82      val mcCabe: Leet[SHA, Bag[Int]] = resources
83          .filter { case (path, _) => path.endsWith(".java") }
84          .map { case (_, resource) => computeMCC(resource) }
85
86      val output: Leet[SHA, Single[Int]] = mcCabe.sum()
87      |
88  }
```

Demo

```
9  ▶  object Sandbox {
10
11      // The foreign function computing MCC on a resource.
12      def computeMCC(resource: Resource): Int = {...}
13
14      // Library imports.
15      import org.toplevel.libs.Natives._
16      import org.toplevel.libs.Gits._
17
18      implicit val engine: Engine = IncrementalParallelEngine.create()
19
20      def main(args: Array[String]): Unit = {
21          val shas: Leet[SHA, Single[SHA]] = git( address = "jwtk/jjwt")
22
23          val resources: Leet[SHA, Bag[(Path, Resource)]] = shas.resources()
24
25          val mcCabe: Leet[SHA, Bag[Int]] = resources
26              .filter { case (path, _) => path.endsWith(".java") }
27              .map { case (_, resource) => computeMCC(resource) }
28
29          val output: Leet[SHA, Single[Int]] = mcCabe.sum()
30
31          output.show()
32      }
33  }
```

Demo

qegal - Google Drive

SANER 20 (+) - Google Slides

Telegram Web

/C:/Data/Repos/topleet/temp/view

/C:/Data/Repos/topleet/t

←

→

↺

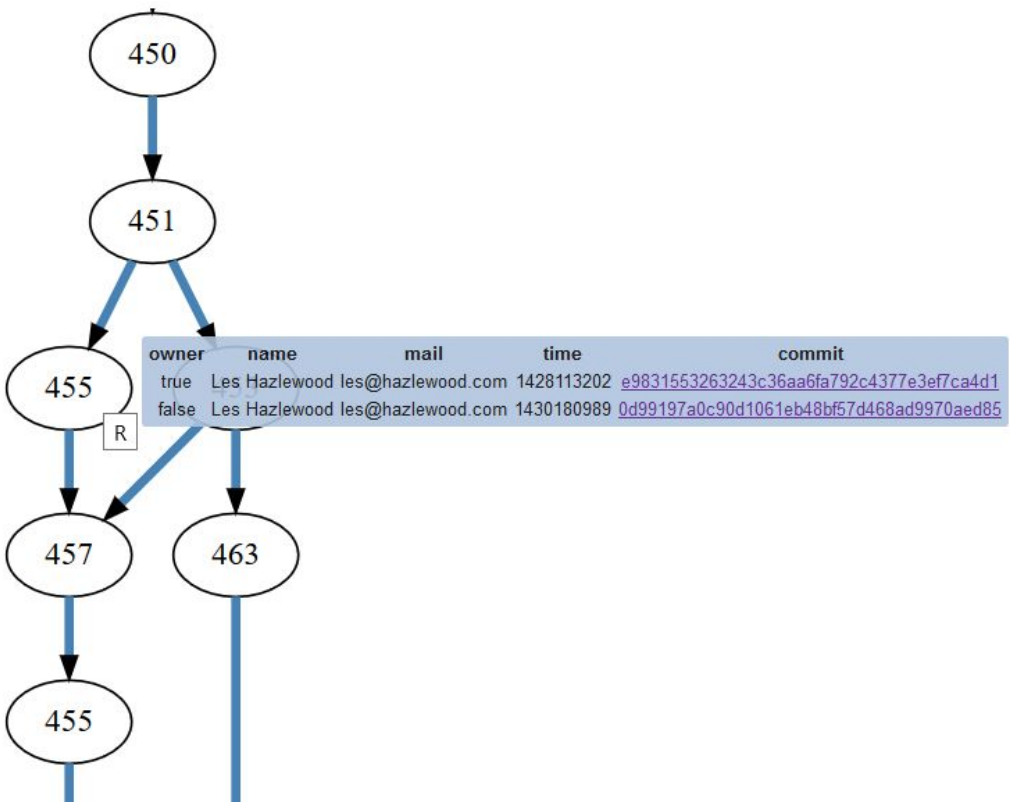
🏠

file:///C:/Data/Repos/topleet/temp/viewer/viewer.html

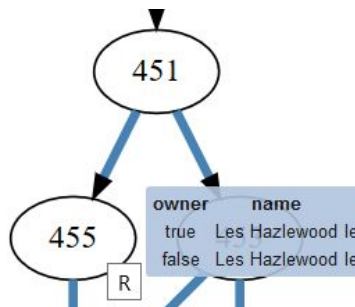
...

🔒

🔍



Demo



```
24 - return DatatypeConverter.printBase64Binary(data);
27 +
28 + if (STANDARD_JVM) {
29 +     return javax.xml.bind.DatatypeConverter.printBase64Binary(data);
30 + }
31 +
32 + if (ANDROID) {
33 +     int flags = android.util.Base64.NO_PADDING | android.util.Base64.NO_WRAP;
34 +     return android.util.Base64.encodeToString(data, flags);
35 + }
36 +
37 + throw new IllegalStateException("Unable to locate a Base64 codec for the current JVM");
25 38 }
26 39
27 40 + @Override
28 41 public byte[] decode(String encoded) {
29 - return DatatypeConverter.parseBase64Binary(encoded);
42 +
43 + if (STANDARD_JVM) {
44 +     return javax.xml.bind.DatatypeConverter.parseBase64Binary(encoded);
45 + }
46 +
47 + if (ANDROID) {
48 +     return android.util.Base64.decode(encoded, android.util.Base64.DEFAULT);
49 + }
50 +
51 + throw new IllegalStateException("Unable to locate a Base64 codec for the current JVM");
30 52 }
31 53
32 54 }
```



Many
Resources on
every Commit



object Sandbox {

// The foreign function computing MCC on a resource.

def computeMCC(resource: Resource): Int = {...}

// Library imports.

import org.topleet.libs.Natives._

import org.topleet.libs.Gits._

implicit val engine: Engine = IncrementalParallelEngine.create()

def main(args: Array[String]): Unit = {

val shas: Leet[SHA, Single[SHA]] = git(address = "jwtk/jjwt")

val resources: Leet[SHA, Bag[(Path, Resource)]] = shas.resources()

val mcCabe: Leet[SHA, Bag[Int]] = resources

.filter { case (path, _) => path.endsWith(".java") }

.map { case (_, resource) => computeMCC(resource) }

val output: Leet[SHA, Single[Int]] = mcCabe.sum()

output.show()

}

Joi

}

}