

Apache AsterixDB

...

Eddie Coberly

Background

- Customer: Apache Software Foundation
 - Technical mentor: Dr. Preston Carman
- Description:
 - Fixed `interval_overlapping()` function for AsterixDB
 - Added functionality for dynamic hints
 - Implemented data deduplication
- Relevance
 - Didn't work before my project

Goals at the beginning

- Project goals:
 - Complete last interval join function
 - Experiment with optimizing different parts of the code
- Criteria for success:
 - Working query
 - Pull request
 - Optimized query (stretch goal)

What was actually accomplished...

- Goals met:
 - Interval join function now deduplicates data before returning it
- Success criteria met:
 - Working query
 - Pull request (in progress)
 - Waiting for confirmation about specific commit process
 - Need affirmative code review rated +2
- Abandoned goals:
 - Optimized query
 - This was a stretch goal anyway

Schedule

| | Projected | Actual |
|--------|---------------------------------------|--|
| Winter | Have PR ready | Create framework for function |
| Spring | Optimize code Update documentation | Complete function Achieve passing tests PR (in progress) |

Tools

- IntelliJ IDEA Java Editor (for coding)
- VMware Workstation 17 (for running AsterixDB on Windows)
 - Abandoned
 - Used Linux lab computer instead
- Gerrit (for submitting changes)
- Jenkins (for automated testing)
 - Not needed

Tools, cont'd

How IntelliJ helped:

- User-friendly
- Lots of keyboard shortcut tips
 - New one every time IDE was opened
- Easy for imports/dependencies
- Code formatting suggestions



Queries: Static vs. Dynamic Hint

QUERY INPUT (5/5)

Default TinyCollege PLAN FORMAT JSON OUTPUT FORMAT JSON QUERY HISTORY use TinyCollege, select f.name as staff, d.name as student from Staff as f, Students...



```
1 use TinyCollege;
2
3 select f.name as staff, d.name as student
4 from Staff as f, Students as d
5 where
6 /*+ range [date("1949-12-31"), date("1974-01-01"), date("2000-01-01"), date("2025-01-01")]*/
7 interval_overlapping(f.employment, d.attendance)
8 order by staff
```

SUCCESS: Execution time: 103.546268ms Elapsed time: 110.985102ms Size: 0.00 Kb

Objects Returned: 0

WARNINGS(0)

CLEAR

EXPLAIN



QUERY OUTPUT



JSON JSONL TABLE TREE PLAN EXPORT

Items per page: 10

0 of 0



0

METADATA INSPECTOR

DATAVERSES

- Default
- Metadata
- TinyCollege

REFRESH

DATASETS

DATATYPES

INDEX

USER DEFINED FUNCTIONS

QUERY INPUT (6/6)

Default

TinyCollege

PLAN FORMAT

JSON

OUTPUT FORMAT

JSON

QUERY HISTORY

use TinyCollege; select f.name as staff, d.name as student from Staff as f, Students... < >

```
1 use TinyCollege;
2
3 select f.name as staff, d.name as student
4 from Staff as f, Students as d
5 where
6 interval_overlapping(f.employment, d.attendance)
7 order by staff
```

ERROR: Code: 1 "java.lang.NullPointerException: Cannot invoke "org.apache.hyracks.dataflow.common.data.partition.range.RangeMap.getSplitCount()" because "this.rangeMap" is null"

WARNINGS(0)

CLEAR

EXPLAIN



METADATA INSPECTOR

DATAVERSES

- Default
- Metadata
- TinyCollege

REFRESH

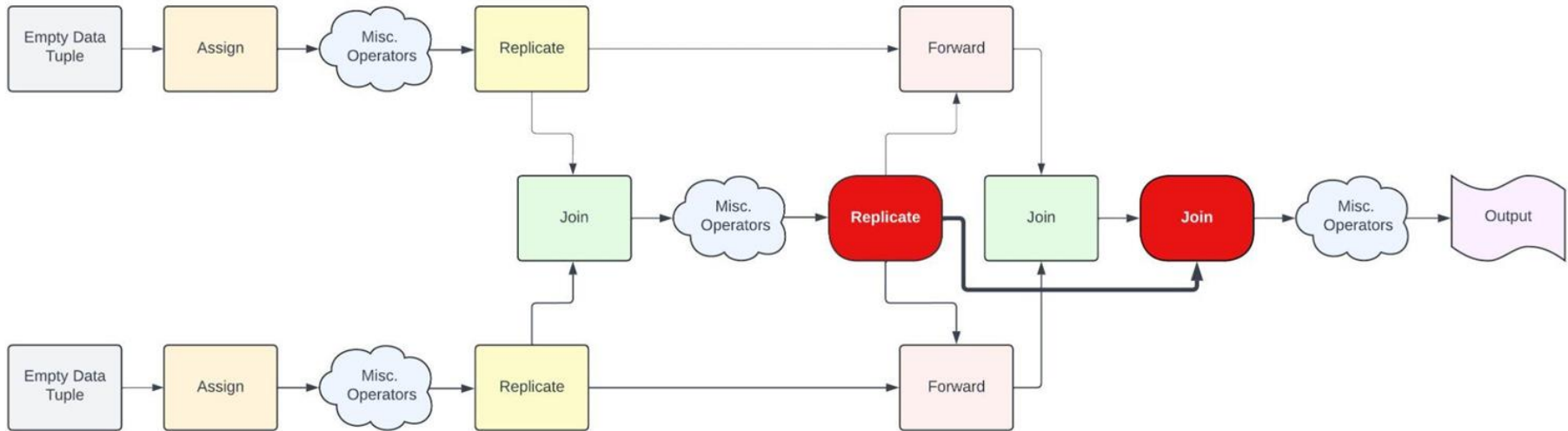
DATASETS

DATATYPES

INDEX

USER DEFINED FUNCTIONS

My solution:



```
93 private final ArrayBackedValueStorage s1 = new ArrayBackedValueStorage();
94
95 2 usages
96 private ISerializerDeserializer<ABoollean> booleanSerde =
97     SerializerDeserializerProvider.INSTANCE.getSerializerDeserializer(BuiltinType.ABOOLEAN);
98
99 ± cobeed
100 @Override
101 public void evaluate(IFrameTupleReference tuple, IPointable result) throws HyracksDataException {
102     resultStorage.reset();
103     if (partition == 0) {
104         booleanSerde.serialize(ABoollean.TRUE, out);
105         result.set(resultStorage);
106         return;
107     }
108     eval0.evaluate(tuple, argPtr0);
109     eval1.evaluate(tuple, argPtr1);
110     eval2.evaluate(tuple, argPtr2);
111
112     if (PointableHelper.checkAndSetMissingOrNull(result, argPtr0, argPtr1, argPtr2)) {
113         return;
114     }
115
116     byte type0 = argPtr0.getTag();
117     byte type1 = argPtr1.getTag();
118
119     if (type0 == ATypeTag.SERIALIZED_INTERVAL_TYPE_TAG && type0 == type1) {
120         argPtr0.getValue(interval0);
121         argPtr1.getValue(interval1);
122         byte intervalType0 = interval0.getType();
123         byte intervalType1 = interval1.getType();
124
125         if (intervalType0 != intervalType1) {
126             throw new IncompatibleTypeException(sourceLoc, getIdentifier(), intervalType0,
127                 intervalType1);
128         }
129     } else if (type0 != ATypeTag.SERIALIZED_INTERVAL_TYPE_TAG) {
130         throw new TypeMismatchException(sourceLoc, getIdentifier(), 0, type0,
131             ATypeTag.SERIALIZED_INTERVAL_TYPE_TAG);
132     } else {
133         throw new IncompatibleTypeException(sourceLoc, getIdentifier(), type0, type1);
134     }
135 }
```

```

472 InnerJoinOperator intervalJoinOp =
473     new InnerJoinOperator(new MutableObject<>(deduplicationTestCondition), leftInputOp, rightInputOp);
474 intervalJoinOp.setSourceLocation(op.getSourceLocation());
475 setSortMergeIntervalJoinOp(intervalJoinOp, fi, sideLeft, sideRight, context, intervalPartitions);
476 intervalJoinOp.setSchema(op.getSchema());
477 context.computeAndSetTypeEnvironmentForOperator(intervalJoinOp);
478
479 Mutable<ILogicalOperator> opRef = new MutableObject<>(op);
480 Mutable<ILogicalOperator> intervalJoinOpRef = new MutableObject<>(intervalJoinOp);
481
482 InnerJoinOperator deduplicationTestOp =
483     new InnerJoinOperator(new MutableObject<>(deduplicationTestCondition), intervalJoinOpRef,
484         exchMAPToDeduplicationJoinOpRef);
485 deduplicationTestOp.setPhysicalOperator(new NestedLoopJoinPOperator(
486     AbstractBinaryJoinOperator.JoinKind.INNER, AbstractJoinPOperator.JoinPartitioningType.BROADCAST));
487 MutableObject<ILogicalOperator> referencePointTestJoinOpRef = new MutableObject<>(deduplicationTestOp);
488 deduplicationTestOp.setSourceLocation(op.getSourceLocation());
489 context.computeAndSetTypeEnvironmentForOperator(deduplicationTestOp);
490 deduplicationTestOp.recomputeSchema();
491 opRef.setValue(referencePointTestJoinOpRef.getValue());
492 op.getInputs().clear();
493 op.getInputs().addAll(deduplicationTestOp.getInputs());
494 op.setPhysicalOperator(deduplicationTestOp.getPhysicalOperator());
495 op.getCondition().setValue(deduplicationTestOp.getCondition().getValue());
496 context.computeAndSetTypeEnvironmentForOperator(op);
497 op.recomputeSchema();
498 }
499

```

1 usage ± cobeed

```

500 @ private static ScalarFunctionCallExpression createIntervalDuplicateTest(AbstractBinaryJoinOperator op,
501     LogicalVariable leftInputVar, LogicalVariable rightInputVar, LogicalVariable rangeMapVar) {
502     // Compute reference interval ID
503     ScalarFunctionCallExpression intervalId = new ScalarFunctionCallExpression(
504         BuiltinFunctions.getBuiltinFunctionInfo(BuiltinFunctions.INTERVAL_JOIN_REMOVE_DUPLICATES),
505         new MutableObject<>(new VariableReferenceExpression(leftInputVar)),
506         new MutableObject<>(new VariableReferenceExpression(rightInputVar)),
507         new MutableObject<>(new VariableReferenceExpression(rangeMapVar)));
508     intervalId.setSourceLocation(op.getSourceLocation());
509
510     return intervalId;
511 }
512

```

```
41 *
42 * T can be of type date, time or datetime.
43 */
7 usages ± cobeed
44 public class ARangeMapSerializerDeserializer implements ISerializerDeserializer<RangeMap> {
45
46     no usages
47     private static final long serialVersionUID = 1L;
48
49     public static final ARangeMapSerializerDeserializer INSTANCE = new ARangeMapSerializerDeserializer();
50
51     1 usage ± cobeed
52     private ARangeMapSerializerDeserializer() {
53     }
54
55     ± cobeed
56     @Override
57     public RangeMap deserialize(DataInput in) throws HyracksDataException {
58         try {
59             int numFields = IntegerSerializerDeserializer.read(in);
60             byte[] bytes = ByteArraySerializerDeserializer.read(in);
61             int[] endOffsets = IntArraySerializerDeserializer.read(in);
62             double[] percentages = DoubleArraySerializerDeserializer.read(in);
63
64             return new RangeMap(numFields, bytes, endOffsets, percentages);
65         } catch (IOException e) {
66             throw HyracksDataException.create(e);
67         }
68     }
69
70     ± cobeed
71     @Override
72     public void serialize(RangeMap instance, DataOutput out) throws HyracksDataException {
73         try {
74             IntegerSerializerDeserializer.write(instance.getNumFields(), out);
75             ByteArraySerializerDeserializer.write(instance.getByteArray(), out);
76             IntArraySerializerDeserializer.write(instance.getEndOffsets(), out);
77             DoubleArraySerializerDeserializer.write(instance.getPercentages(), out);
78         } catch (IOException e) {
79             throw HyracksDataException.create(e);
80         }
81     }
82 }
```

33

± cobeed

34 public class RangeMapSerializerDeserializerTest {

35

± cobeed

@Test

36

37 public void test() throws HyracksDataException {

38

39 ArrayBackedValueStorage abvs = new ArrayBackedValueStorage();

40 DataOutput out = abvs.getDataOutput();

41

42 AInt64SerializerDeserializer intSerde = AInt64SerializerDeserializer.INSTANCE;

43 int[] offsets = new int[3];

44 double[] percentages = new double[3];

45 try {

46 // Loop over list of integers

47 for (int i = 0; i < 3; i++) {

48 intSerde.serialize(new AMutableInt64(value: (long) i + 10), out);

49 offsets[i] = abvs.getLength();

50 percentages[i] = i * 1.3;

51 }

52 } catch (IOException e) {

53 throw HyracksDataException.create(e);

54 }

55 RangeMap rangeMap = new RangeMap(numFields: 1, abvs.getByteArray(), offsets, percentages);

56

57 ArrayBackedValueStorage abvsResult = new ArrayBackedValueStorage();

58 DataOutput outResult = abvsResult.getDataOutput();

59

60 ARangeMapSerializerDeserializer rangeMapSerde = ARangeMapSerializerDeserializer.INSTANCE;

61 rangeMapSerde.serialize(rangeMap, outResult);

62

63 DataInput dataIn = new DataInputStream(new ByteArrayInputStream(abvsResult.getByteArray(), offset: 0, abvsResult.getLength()));

64 // check output

65 RangeMap rangeMapResult = rangeMapSerde.deserialize(dataIn);

66 Assert.assertEquals(rangeMapResult.getNumFields(), rangeMap.getNumFields());

67 Assert.assertArrayEquals(rangeMapResult.getByteArray(), rangeMap.getByteArray());

68 Assert.assertArrayEquals(rangeMapResult.getEndOffsets(), rangeMap.getEndOffsets());

69 Assert.assertArrayEquals(rangeMapResult.getPercentages(), rangeMap.getPercentages(), delta: 0.01);

70 Assert.assertEquals(rangeMapResult, rangeMap);

71 }

72 }

And then...

QUERY INPUT (1/1)

Default

PLAN FORMAT

OUTPUT FORMAT

QUERY HISTORY

TinyCollege

JSON

JSON

use TinyCollege; select f.name as staff, d.name as student from Staff as f, Students...



```
1 use TinyCollege;
2
3 select f.name as staff, d.name as student
4 from Staff as f, Students as d
5 where
6 interval_overlapping(f.employment, d.attendance)
7 order by staff
```

SUCCESS: Execution time: 517.933938ms Elapsed time: 576.303264ms Size: 175.84 Kb

Objects Returned: 4617

WARNINGS(0)

CLEAR

EXPLAIN



QUERY OUTPUT



JSON

JSONL

TABLE

TREE

PLAN

EXPORT

Items per page: 10

1 - 10 of 4617



```
{"staff":"AEJO","student":"jfbj"}
{"staff":"AEJO","student":"jybx"}
{"staff":"AEJO","student":"qwpa"}
{"staff":"AEJO","student":"krgo"}
{"staff":"AEJO","student":"onya"}
{"staff":"AEJO","student":"egae"}
{"staff":"AEJO","student":"lntr"}
{"staff":"AEJO","student":"mfmo"}
{"staff":"AEJO","student":"ybsl"}
```

Challenges – Going In

- Expected:
 - Installing and configuring IDE and other tools
 - Reading and understanding legacy code
 - Learning Java and navigating the IDE
 - Tests run for a long time
- Unexpected:
 - Failed tests
 - Running software on Windows
 - Understanding architecture

Challenges – Going Out

- Expected:
 - Installing AsterixDB
 - Reading and understanding legacy code
 - Time-consuming build process
 - Learning Java and navigating the IDE
- Unexpected
 - Meeting with project mentor
 - Once or twice a week usually, but we had plenty of schedule conflicts
 - Confined to the CS lab
 - Staying awake
 - Keeping up with the hours
 - Submitting code
 - Documentation is sorely out-of-date

Capstone Experience

- Learn a new programming language
 - Java
- Use specialized IDE and suite of developer tools
 - Debugging
- Understand and contribute to legacy code
- Use existing framework to accomplish the goal
- Contact developers
- Regularly report to a customer

Resources

- <https://asterixdb.apache.org/dev-setup.html>
- <https://asterixdb.apache.org/pushing.html>
- <https://asterixdb.apache.org/docs/0.9.9/sqlpp/builtins.html>
- https://asterixdb.apache.org/docs/0.9.9/interval_join.html

Questions?