

Week6 WAL / Manifest

Made by Suhwan Shin, Isu Kim, Seyeon Park

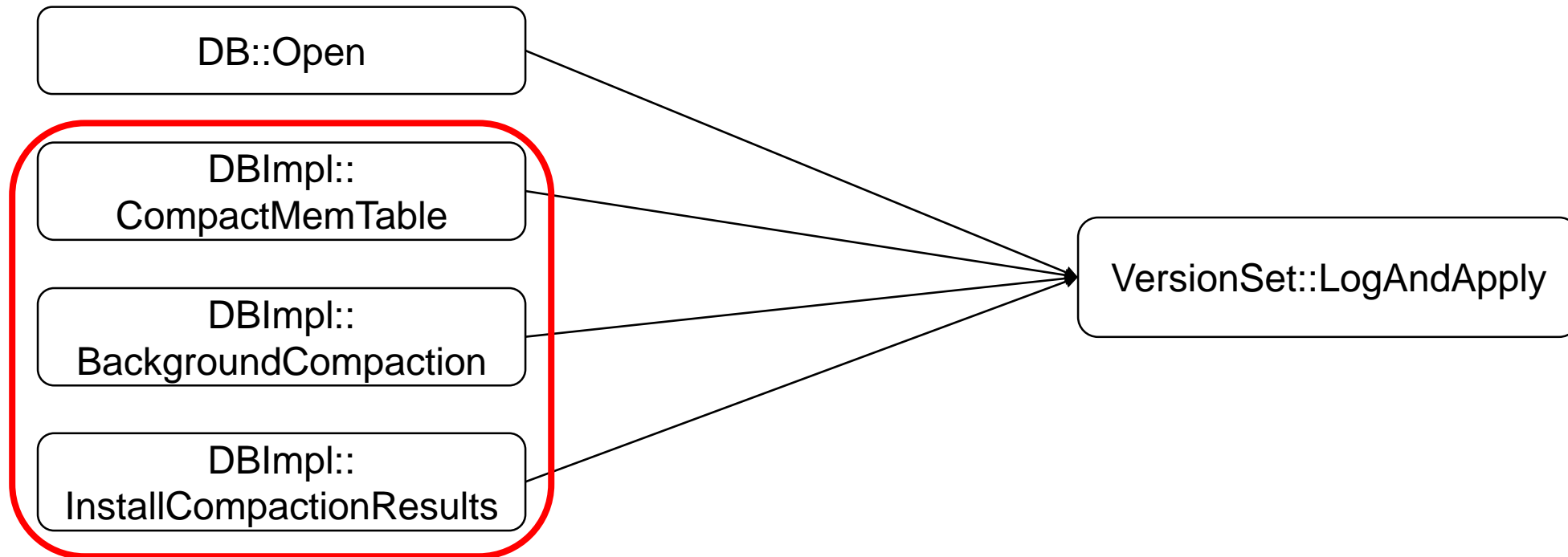
E-Mail: tlstnghks77@dankook.ac.kr

Content

- Code flow
 - VersionSet
- Code flows & Example
 - `log::Writer::AddRecord`
 - `log::Writer::EmitPhysicalRecord`
 - `PosixWritableFile::Append`
 - Put example
- Log format
- Code flow
 - `log::reader`

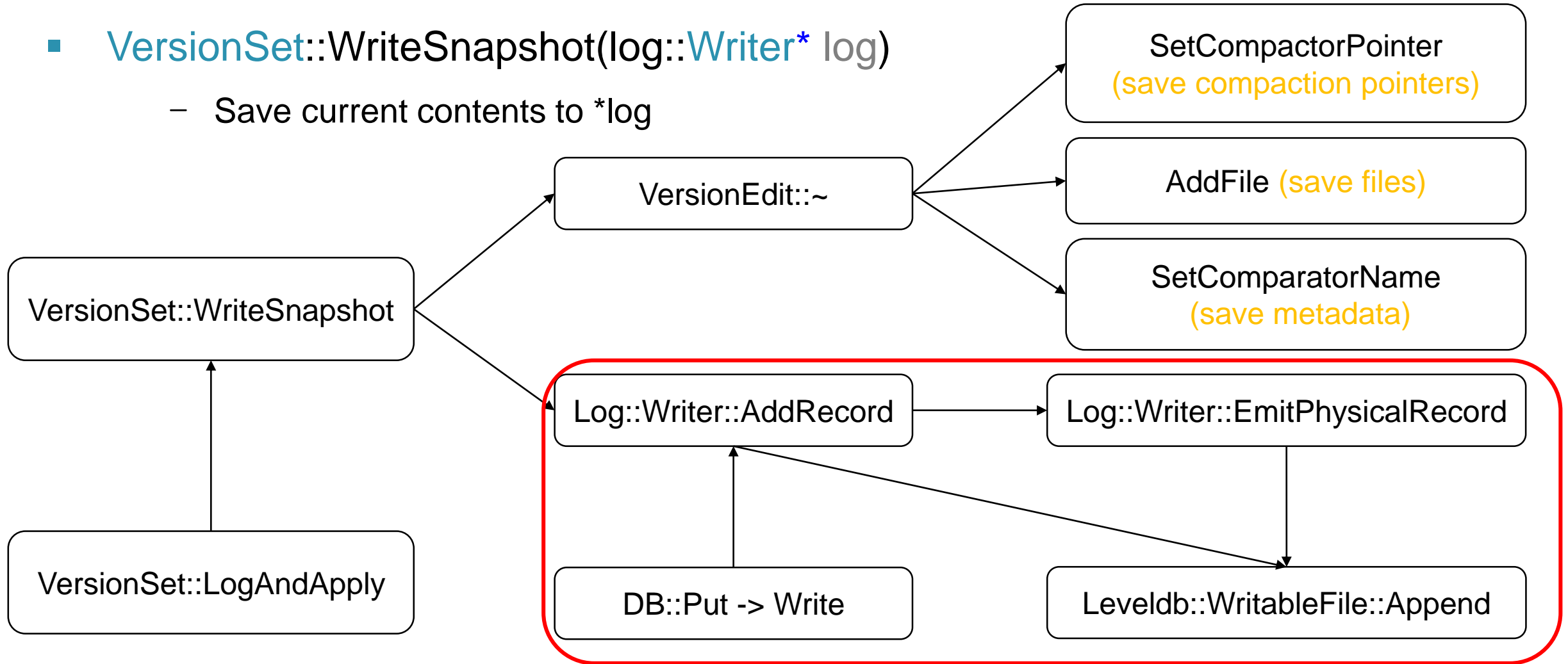
VersionSet::LogAndApply

- Status VersionSet::LogAndApply(VersionEdit* edit, port::Mutex* mu)
 - Apply *edit to the current version to form a new descriptor that is both saved to persistent state and installed as the new current version.



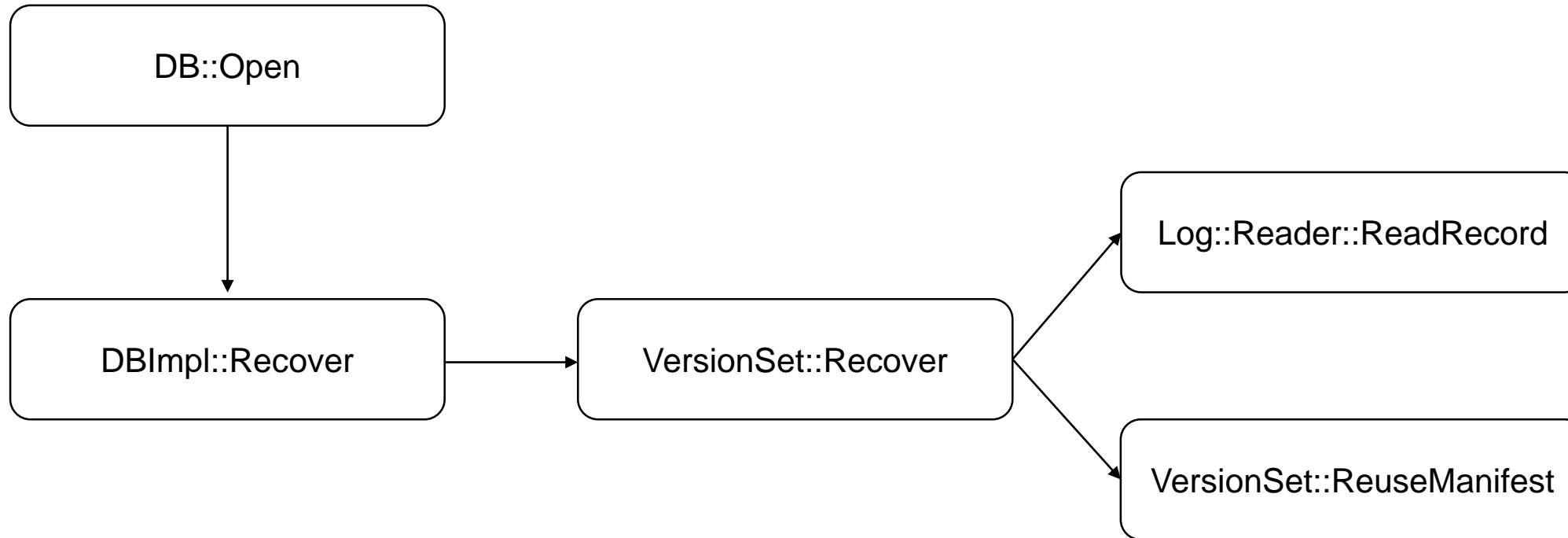
VersionSet::WriteSnapshot

- VersionSet::WriteSnapshot(log::Writer* log)
 - Save current contents to *log



VersionSet::Recover

- `VersionSet::Recover(bool* save_manifest)`
 - Recover the last saved descriptor from persistent storage.

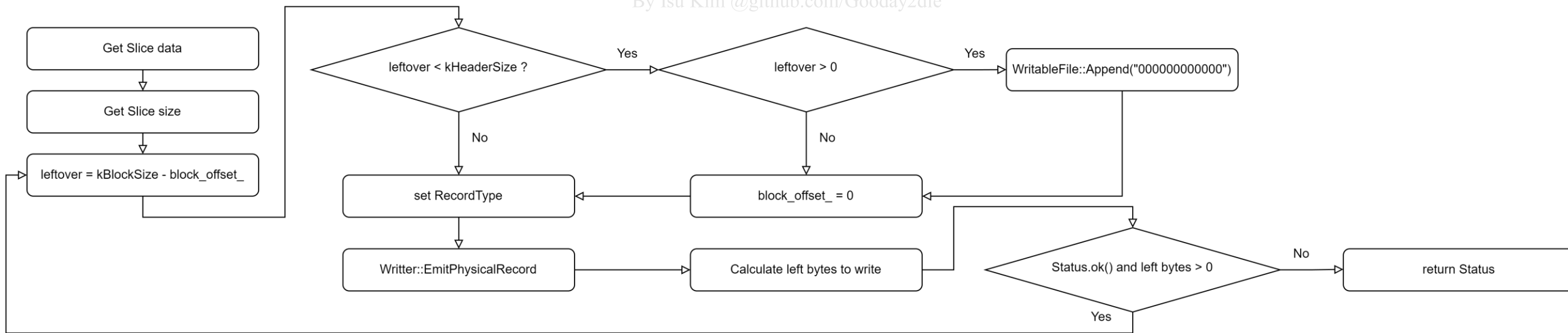


log::Writer::AddRecord

- Iterate and write Slice data till Slice is empty
- Defined in log_writer.cc @ line 34
- Checks Record Type and writes data using EmitPhysicalRecord.

<https://github.com/DKU-StarLab/leveldb-study/>

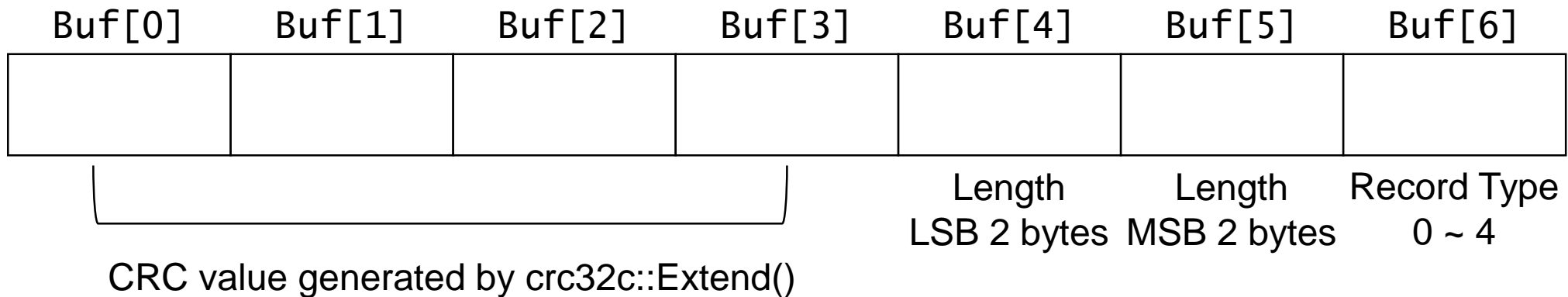
By Isu Kim @github.com/Gooday2die



log::Writer::EmitPhysicalRecord

- Generates Header in following format
- Defined in log_writer.cc @ line 86
- kHeaderSize defaults to 7
- Will be written using `PosixWritableFile::Append`

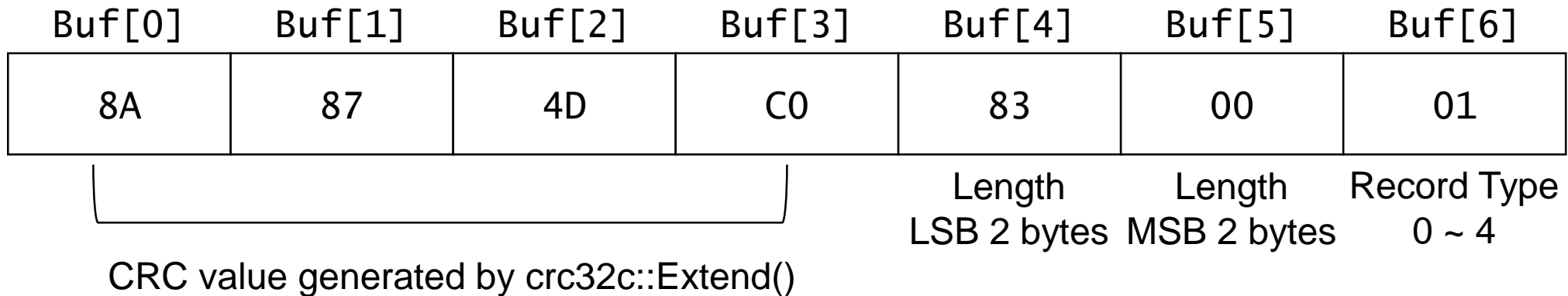
<https://github.com/DankookUniversity/ereldp-study>
By Isu Kim @github.com/Goody2die



log::Writer::EmitPhysicalRecord - Example

- .log file generated by db_bench
- **CRC Value** : 0xC04D878A
- **Length** : 0x0083 (131 bytes)
- **Record Type** : 0x01 (kFullType)

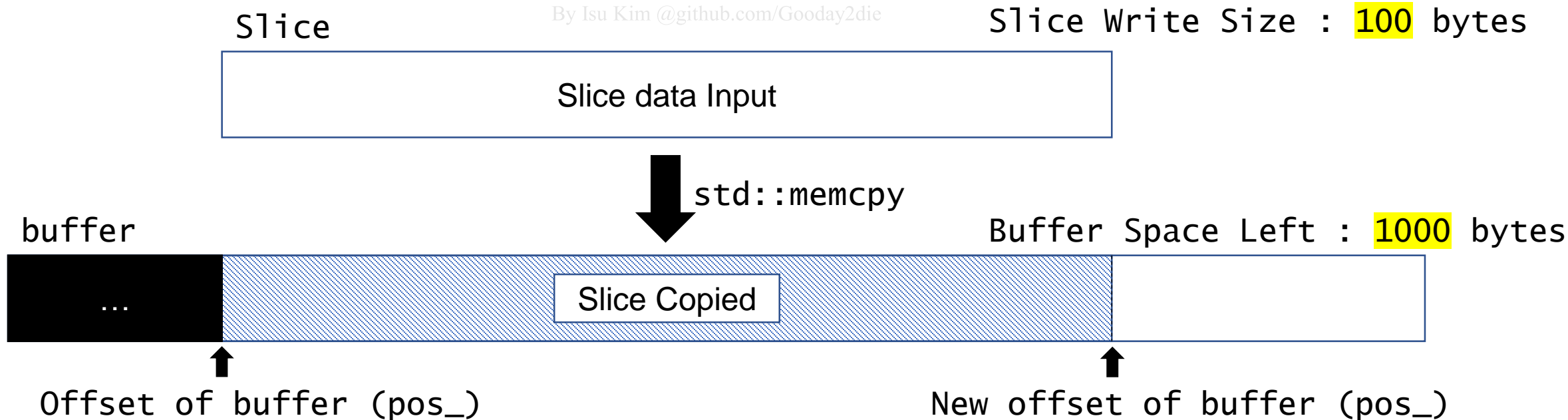
<https://github.com/DKU-StarLab/leveldb-study/>
By Isu Kim @github.com/Gooday2die



PosixWritableFile::Append – Case 1

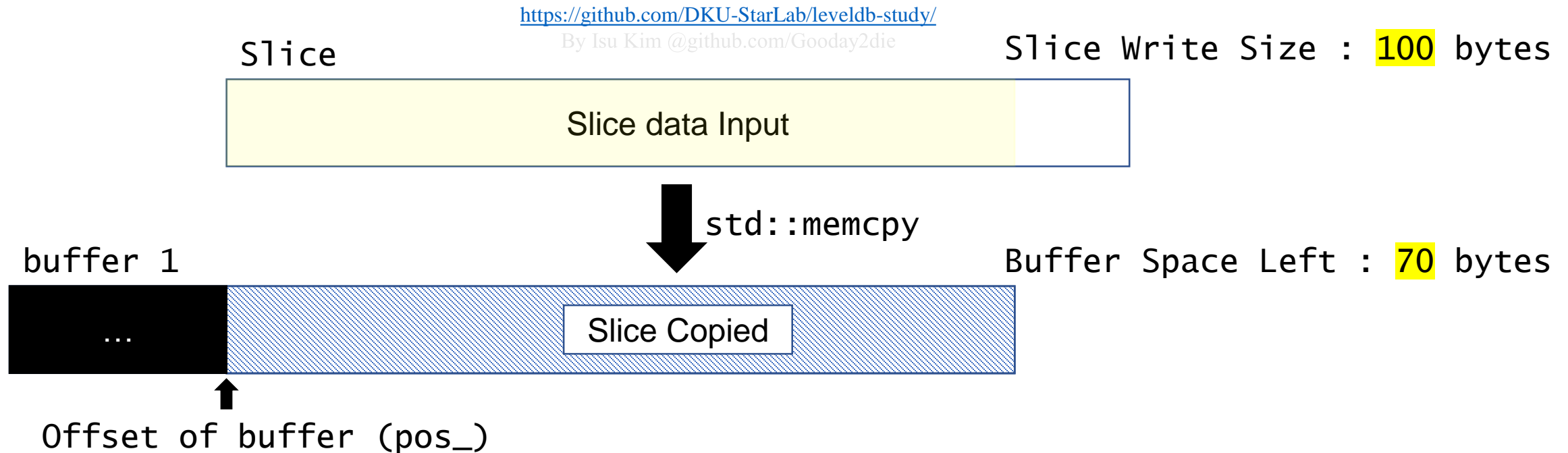
- When buffer size of WritableFile was **bigger** than write size.
 - Copy all Slice data.
 - Calculate buffer offset(pos_) and return Status::OK().

<https://github.com/DKU-StarLab/leveldb-study/>
By Isu Kim @github.com/Gooday2die



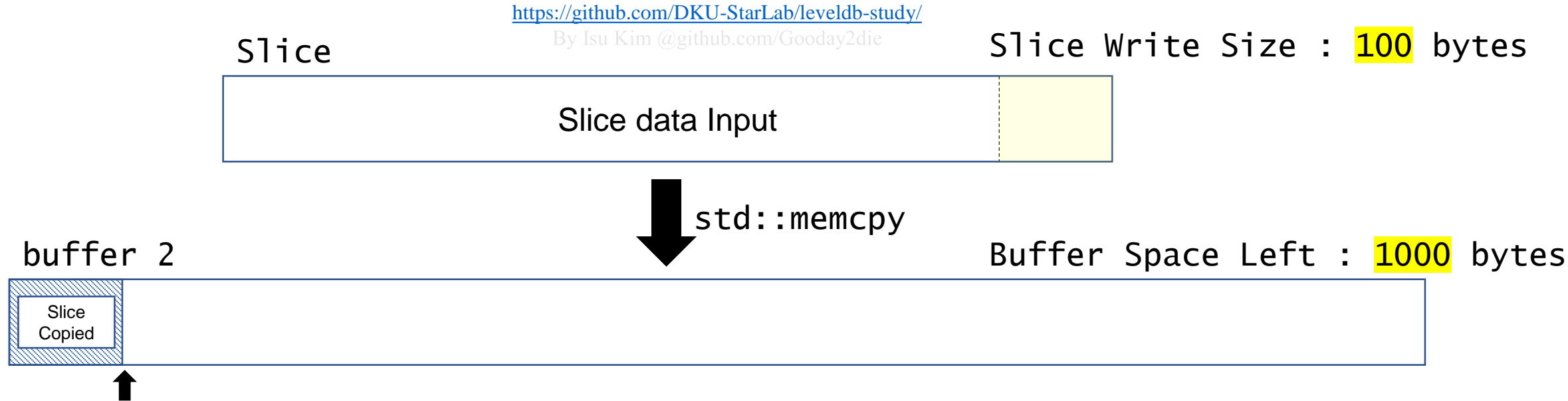
PosixWritableFile::Append – Case 2.1

- When buffer size of WritableFile was **smaller** than write size.
 - Copy as much data as possible from Slice.
 - Perform FlushBuffer()



PosixWritableFile::Append – Case 2.1

- When buffer size of WritableFile was **bigger** than write size.
 - Copy all Slice data.
 - Calculate buffer offset(pos_) and return Status::OK().

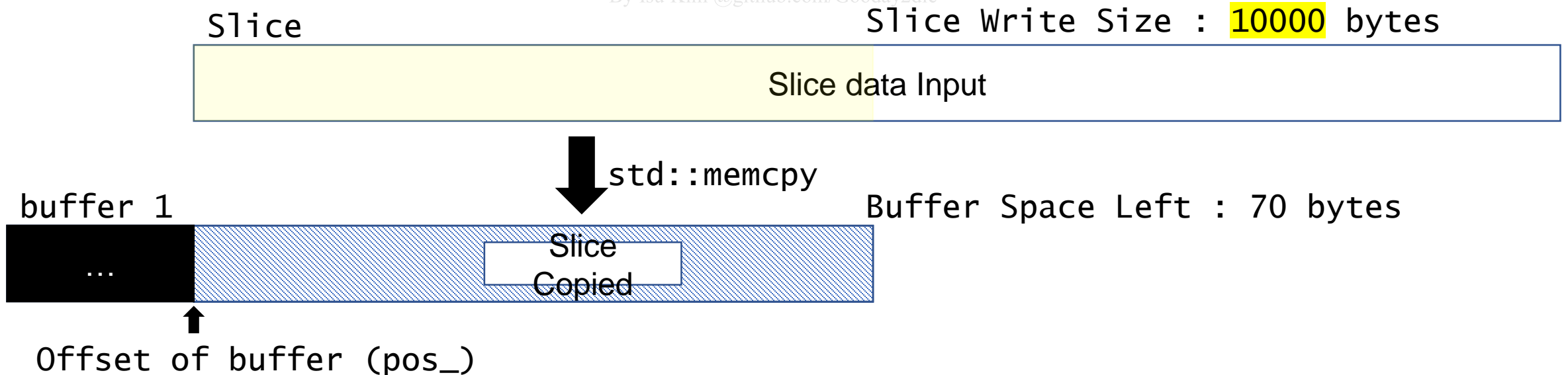


PosixWritableFile::Append – Case 2.2

- When buffer size of WritableFile was **smaller** than write size.
 - Copy as much data as possible from Slice.
 - Perform FlushBuffer()

<https://github.com/DKU-StarLab/leveldb-study/>

By Isu Kim @github.com/Goody2die

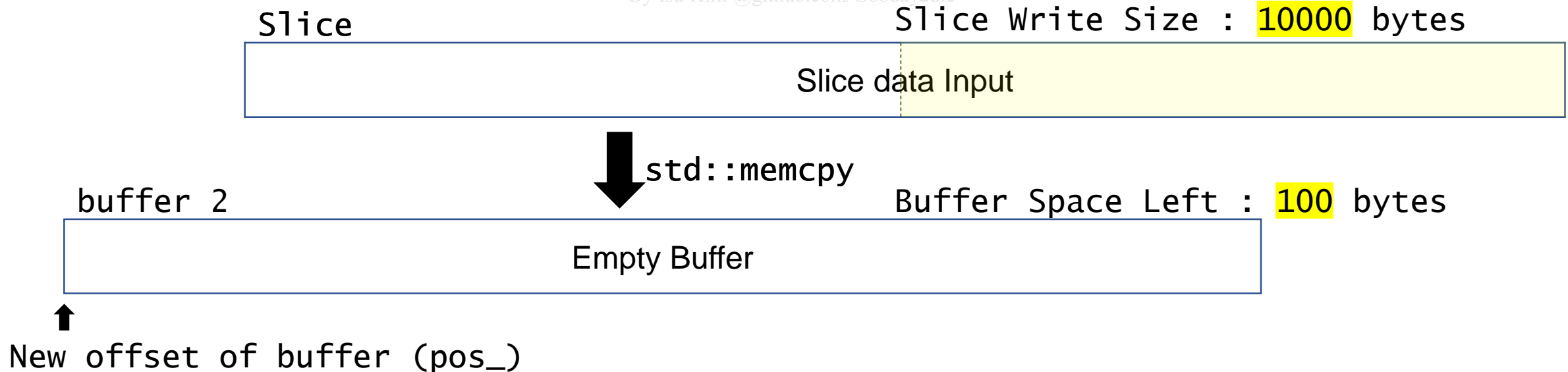


PosixWritableFile::Append – Case 2.2

- When buffer size of WritableFile was **bigger** than write size.
 - Write left data using WriteUnbuffered();
 - Return Status::OK().

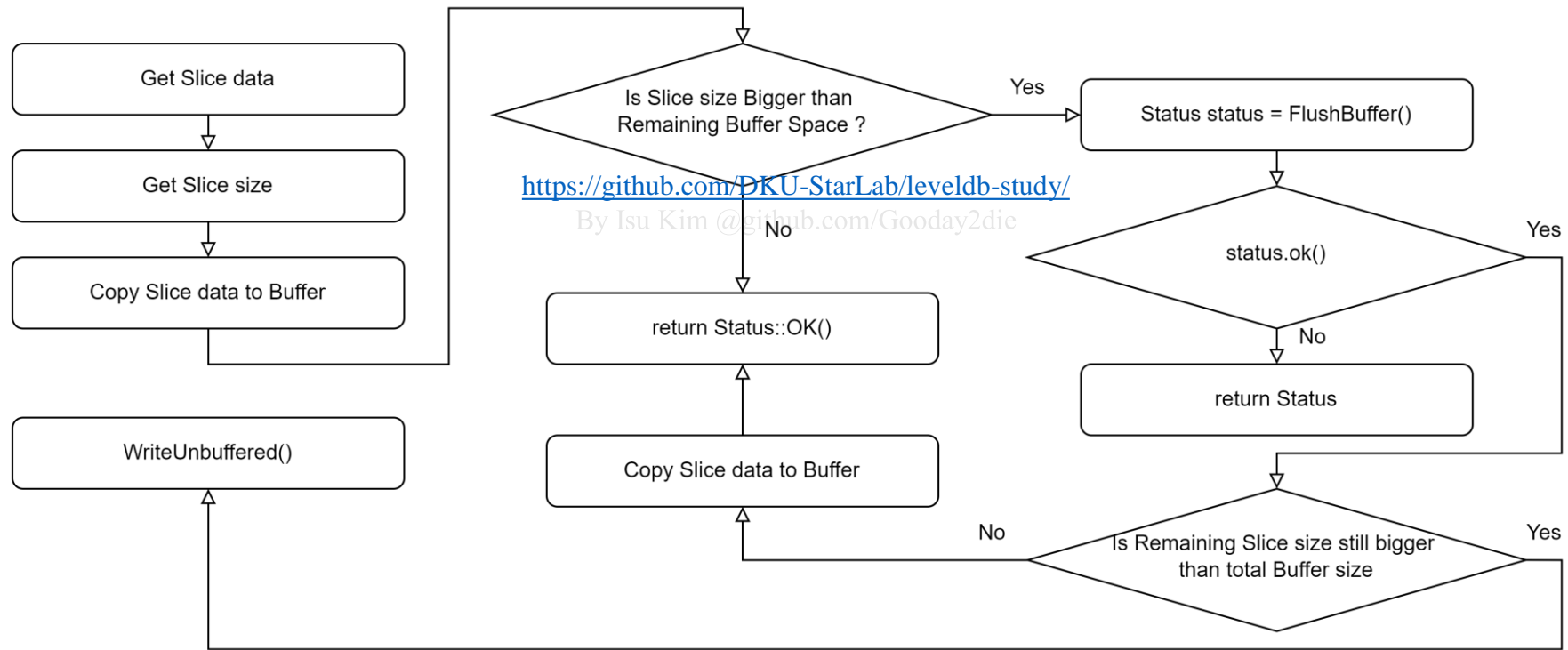
<https://github.com/DKU-StarLab/leveldb-study/>

By Isu Kim @github.com/Gooday2die



PosixWritableFile::Append - Flowchart

- Flowchart of PosixWritableFile::Append



WAL Example for Put - Header

- Tested 'put' using {"A": "Hello world!", "B": "Good bye world!", "C": "I am hungry"}

```

0000h: C0 55 16 4E 1C 00 01 01 00 00 00 00 00 00 00 01 AU.N.....
0010h: 00 00 00 01 01 41 0C 48 65 6C 6C 6F 20 77 6F 72 ....A.Hello wor
0020h: 6C 64 21 8E C7 D9 9D 1F 00 01 02 00 00 00 00 00 ld!ŽÇÜ.....
0030h: 00 00 01 00 00 00 01 01 42 0F 47 6F 6F 64 20 62 .....B.Good b
0040h: 79 65 20 77 6F 72 6C 64 21 26 5D 13 F3 1B 00 01 ye world!&].ó...
0050h: 03 00 00 00 00 00 00 00 01 00 00 00 00 00 00 00 http://github.com/DK-STARLab/leveldb-study/
0060h: 49 20 61 6D 20 68 75 6E 67 72 79 79 79 79 79 79 By Isu Kim @github.com I am hungry
    
```

.log file

```

int main (void) {
    dbTest* db = new dbTest();

    db->putValue("A", "Hello world!");
    db->putValue("B", "Good bye world!");
    db->putValue("C", "I am hungry");

    delete(db);
    return 0;
}
    
```

- Example with {"A": "Hello World!"}

Buf[0]	Buf[1]	Buf[2]	Buf[3]	Buf[4]	Buf[5]	Buf[6]
C0	55	16	4E	1C	00	01
				Length	Length	Record Type
				LSB 2 bytes	MSB 2 bytes	0 ~ 4

CRC value generated by crc32c::Extend()

→ Checksum : 0x4E1655C0 / Length: 0x1C(28) bytes / Record Type : kFullType

WAL Example for Put – Entry Number

0000h:	C0 55 16 4E	1C 00 01	01 00	00 00 00	00 00 00 01	ÀU.N.....
0010h:	00 00 00 01	01 41 0C 48	65 6C 6C 6F	20 77 6F 72A.Hello wor	
0020h:	6C 64 21 8E	C7 D9 9D 1F	00 01	02 00	00 00 00 00	ld!ŽÇÛ.....
0030h:	00 00 01 00	00 00 01 01	42 0F 47 6F	6F 64 20 62B.Good b	
0040h:	79 65 20 77	6F 72 6C 64	21 26 5D 13	F3 1B 00 01	ye world!&].ó...	
0050h:	03 00	00 00 00 00	01 00 00 00	00 00 00 00C.	
0060h:	49 20 61 6D	20 68 75 6E	67 72 79	00 00 00 00	I am hungry	

.log file

```
int main (void) {  
    dbTest* db = new dbTest();  
  
    db->putValue("A", "Hello world!");  
    db->putValue("B", "Good bye world!");  
    db->putValue("C", "I am hungry");  
  
    delete(db);  
    return 0;  
}
```

- For “Hello World!” the Entry number was 0x0001
- For “Good bye world!” the Entry number was 0x0002
- For “I am hungry” the Entry number was 0x0003

WAL Example for Put – Mystery

```
0000h: C0 55 16 4E 1C 00 01 01 00 00 00 00 00 00 01 ÀU.N.....
0010h: 00 00 00 01 01 41 0C 48 65 6C 6C 6F 20 77 6F 72 .....A.Hello wor
0020h: 6C 64 21 8E C7 D9 9D 1F 00 01 02 00 00 00 00 ld!ŽÙ.....
0030h: 00 00 01 00 00 00 01 01 42 0F 47 6F 6F 64 20 62 .....B.Good b
0040h: 79 65 20 77 6F 72 6C 64 21 26 5D 13 F3 1B 00 01 ye world!&].ó...
0050h: 03 00 00 00 00 00 00 00 01 00 00 00 00 00 00 C.
0060h: 49 20 61 6D 20 68 75 6E 67 72 79 20 20 20 20 I am hungry
```

.log file

```
int main (void) {  
    dbTest* db = new dbTest();  
  
    db->putValue("A", "Hello world!");  
    db->putValue("B", "Good bye world!");  
    db->putValue("C", "I am hungry");  
  
    delete(db);  
    return 0;  
}
```

- Seems like batch information.
- But needs confirmation.

WAL Example file for Put - Data

- Tested 'put' for {"A": "Hello world!"}

```
0000h: C0 55 16 4E 1C 00 01 01 00 00 00 00 00 00 00 01 AU.N.....
0010h: 00 00 00 01 01 41 0C 48 65 6C 6C 6F 20 77 6F 72 .....A.Hello wor
0020h: 6C 64 21 8E C7 D9 9D 1F 00 01 02 00 00 00 00 00 ld!ŽÇÛ.....
0030h: 00 00 01 00 00 00 01 01 42 0F 47 6F 6F 64 20 62 .....B.Good b
0040h: 79 65 20 77 6F 72 6C 64 21 26 5D 13 F3 1B 00 01 ye world!&].ó...
0050h: 03 00 00 00 00 00 00 00 01 00 00 00 00 00 00 00 http://github.com/DK-STARLab/leveldb-study/
0060h: 49 20 61 6D 20 68 75 6E 67 72 79 I am hungry
```

.log file

```
int main (void) {
    dbTest* db = new dbTest();

    db->putValue("A", "Hello world!");
    db->putValue("B", "Good bye world!");
    db->putValue("C", "I am hungry");

    delete(db);
    return 0;
}
```

41	0C	48	65	6C	6C	6F	20	77	6F	72	6C	64	21
----	----	----	----	----	----	----	----	----	----	----	----	----	----

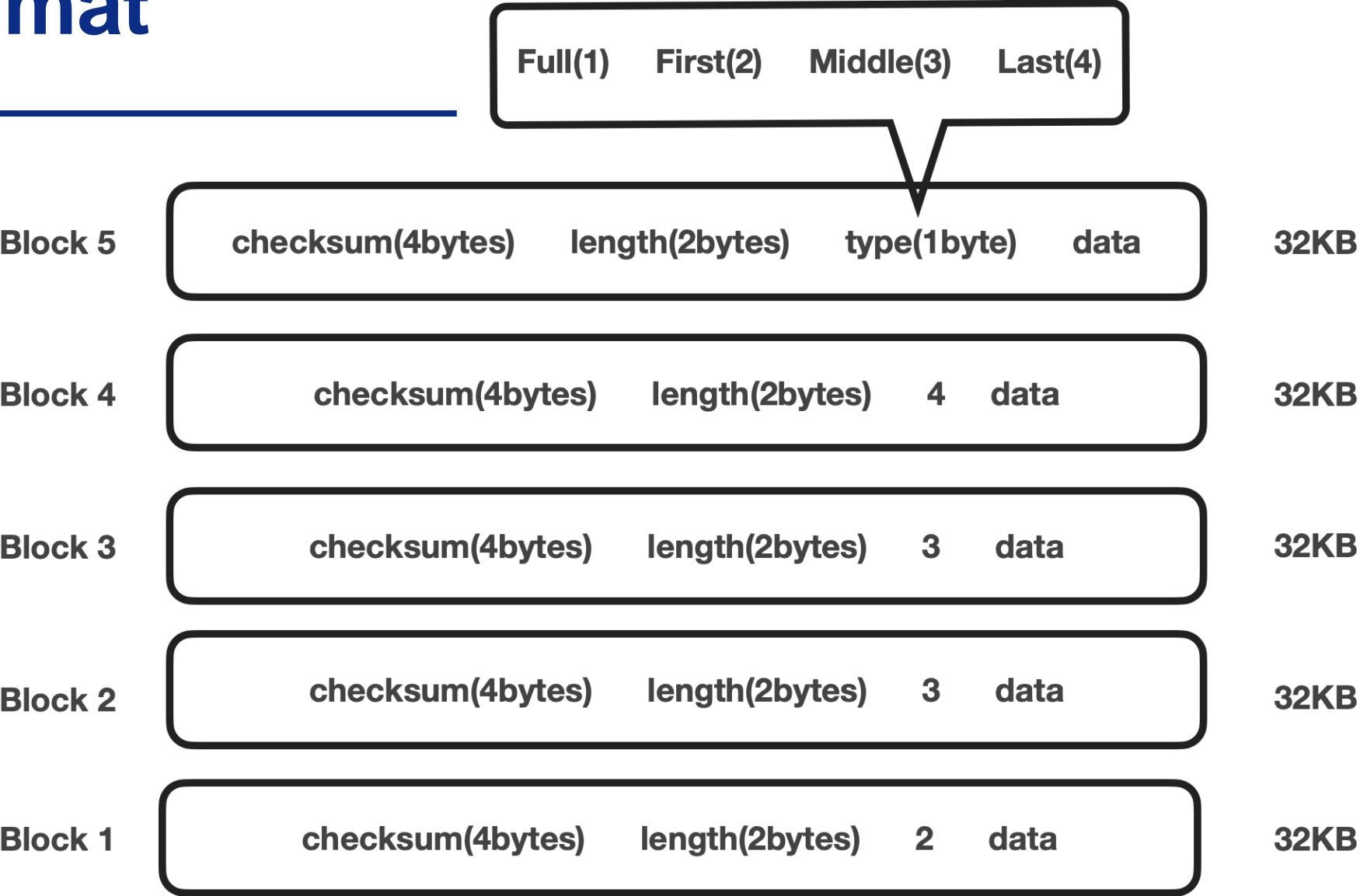


A	^L	H	e	l	l	o		w	o	r	l	d	!
---	----	---	---	---	---	---	--	---	---	---	---	---	---



Form Feed

Log format



⋮

Code flow - log::reader

Class Reader

Class Reporter

Reader(SequentialFile* file, Reporter* reporter, **bool** checksum,
uint64_t initial_offset)

: file_(file),
reporter_(reporter),
checksum_(checksum),
backing_store_(**new char**[kBlockSize]),
buffer_(),
eof_(**false**),
last_record_offset_(0),
end_of_buffer_offset_(0),
initial_offset_(initial_offset),
resyncing_(initial_offset > 0) {}

functions:

bool **Reader::ReadRecord**(Slice* record, std::string* scratch)

bool **Reader::SkipToInitialBlock**()

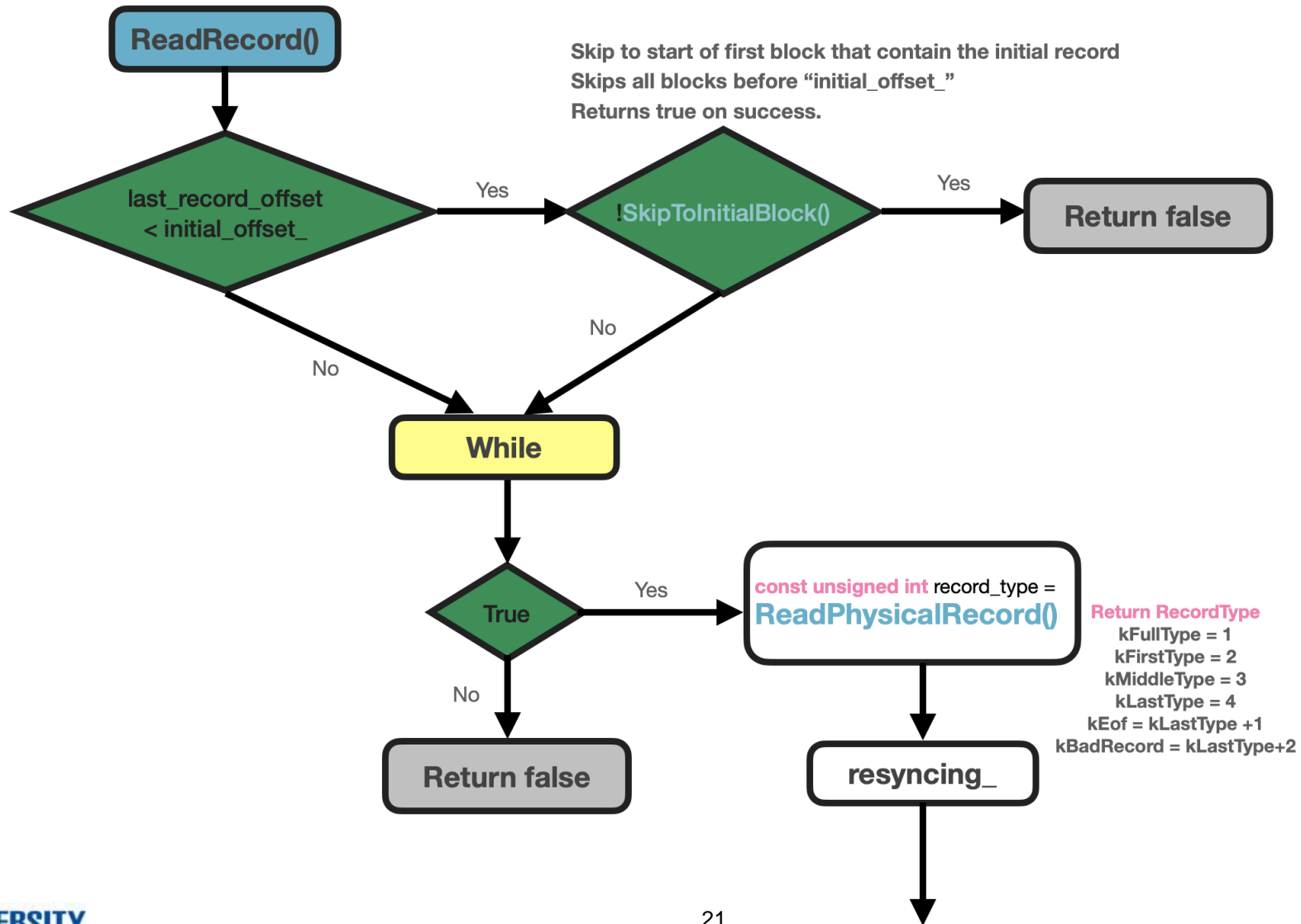
unsigned int **Reader::ReadPhysicalRecord**(Slice* result)

uint64_t **Reader::LastRecordOffset**() { **return** last_record_offset_; }

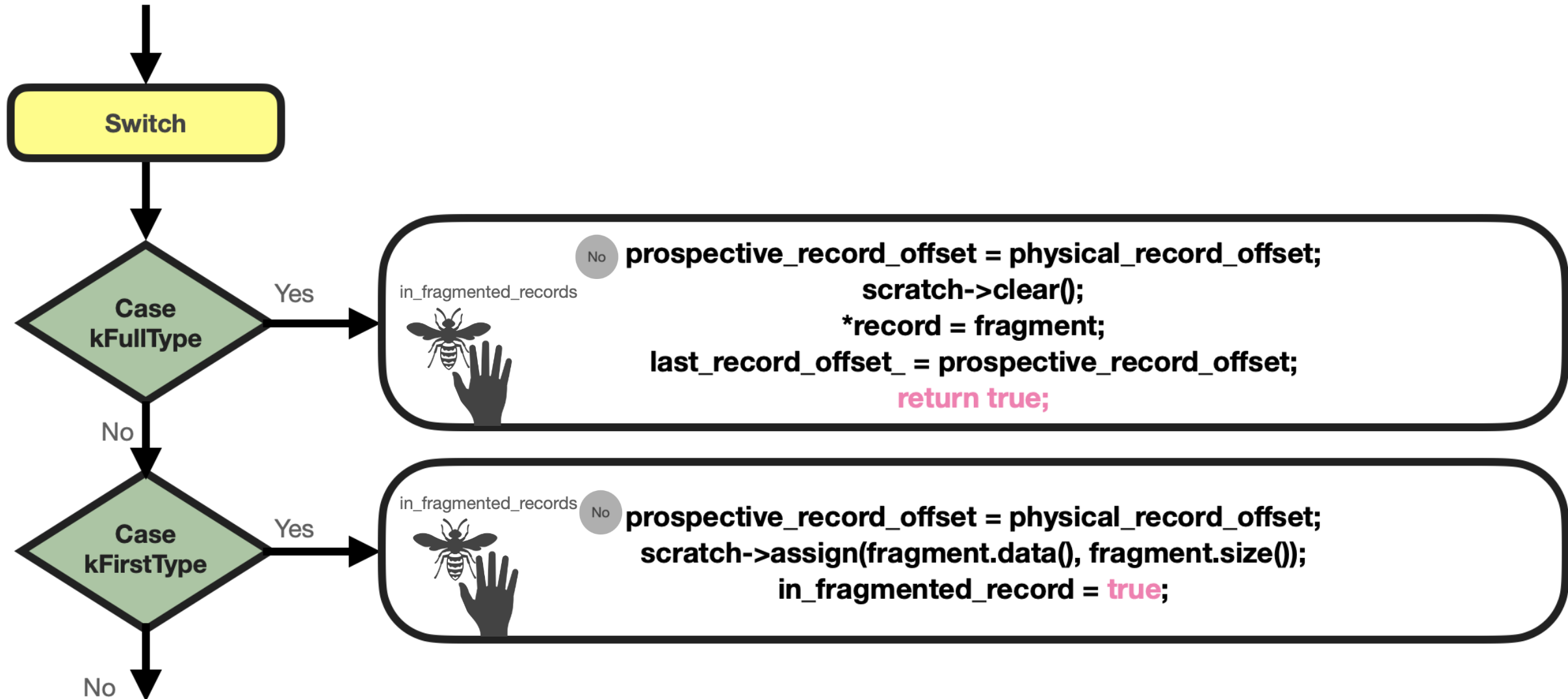
void **Reader::ReportDrop**(uint64_t bytes, **const** Status& reason)

void **Reader::ReportCorruption**(uint64_t bytes, **const char*** reason)

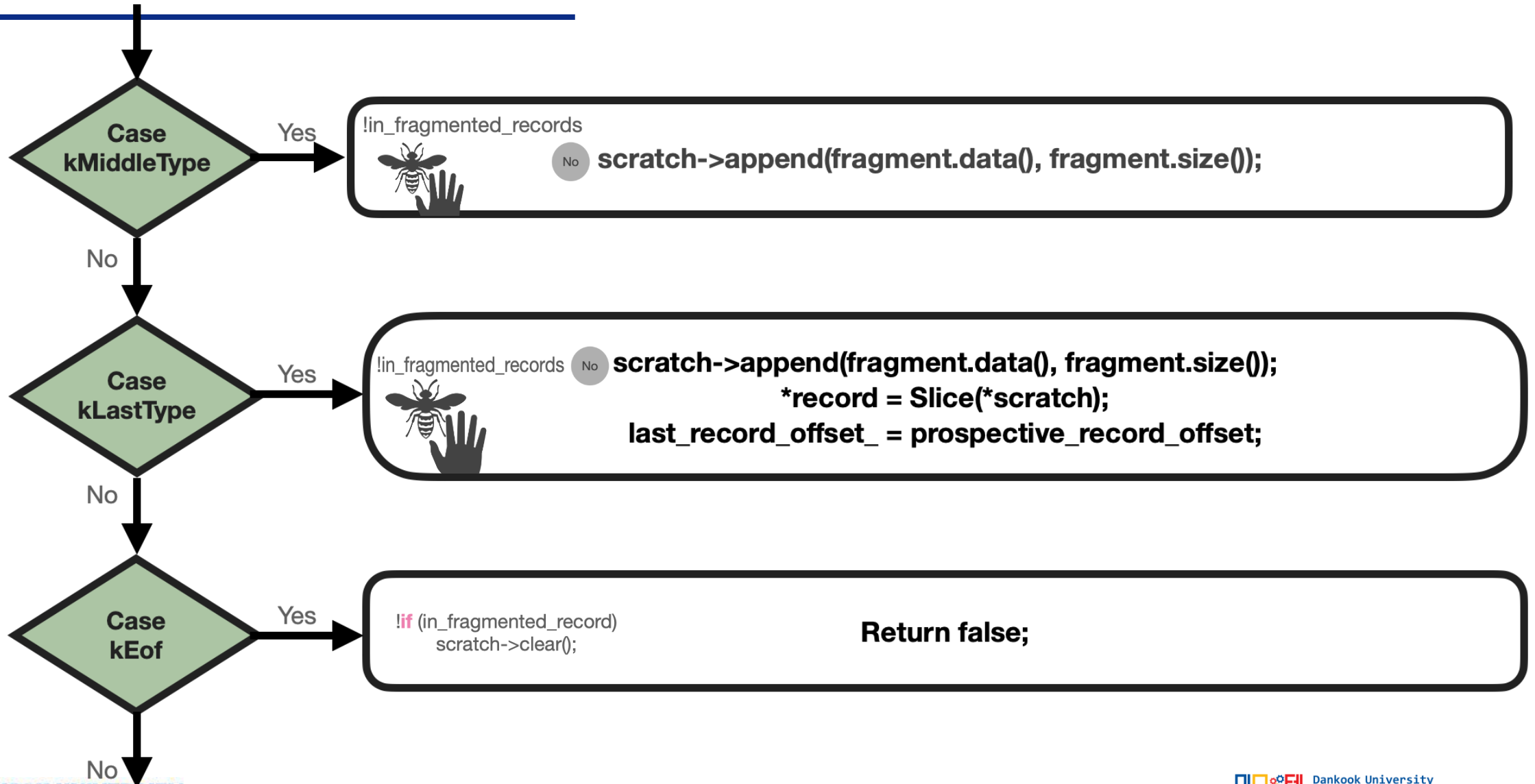
Code flow - log::reader



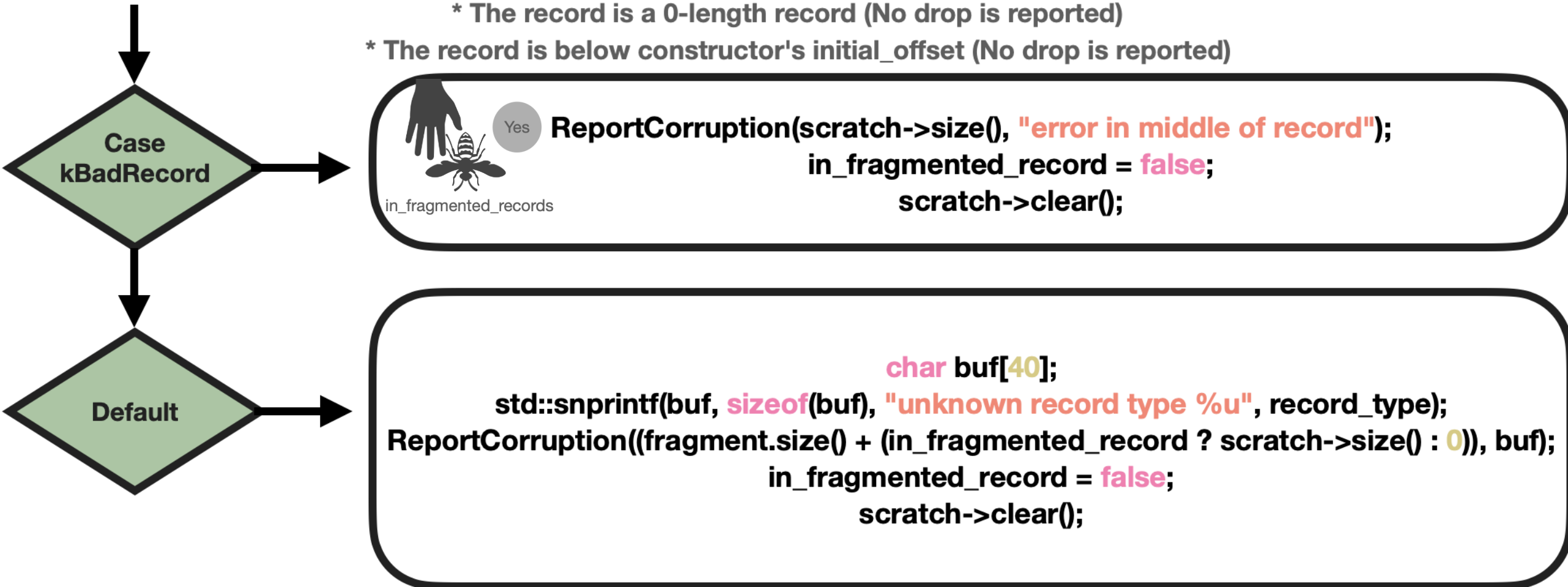
Code flow - log::reader



Code flow - log::reader



Code flow - log::reader



Question

