

---

# Table of Contents

Preface Or: What Are You Getting Yourself Into Here?.....	vii
---	-----

---

## Part I. The Beam Model

<b>1. Streaming 101.....</b>	<b>3</b>
Terminology: What Is Streaming?	4
On the Greatly Exaggerated Limitations of Streaming	6
Event Time Versus Processing Time	9
Data Processing Patterns	12
Bounded Data	12
Unbounded Data: Batch	13
Unbounded Data: Streaming	14
Summary	22
<b>2. The <i>What</i>, <i>Where</i>, <i>When</i>, and <i>How</i> of Data Processing.....</b>	<b>25</b>
Roadmap	26
Batch Foundations: <i>What</i> and <i>Where</i>	28
<i>What</i> : Transformations	28
<i>Where</i> : Windowing	32
Going Streaming: <i>When</i> and <i>How</i>	34
<i>When</i> : The Wonderful Thing About Triggers Is Triggers Are Wonderful Things!	34
<i>When</i> : Watermarks	39
<i>When</i> : Early/On-Time/Late Triggers FTW!	44
<i>When</i> : Allowed Lateness (i.e., Garbage Collection)	47
<i>How</i> : Accumulation	51
Summary	55

<b>3. Watermarks.....</b>	<b>59</b>
Definition	59
Source Watermark Creation	62
Perfect Watermark Creation	64
Heuristic Watermark Creation	65
Watermark Propagation	67
Understanding Watermark Propagation	69
Watermark Propagation and Output Timestamps	75
The Tricky Case of Overlapping Windows	80
Percentile Watermarks	81
Processing-Time Watermarks	84
Case Studies	86
Case Study: Watermarks in Google Cloud Dataflow	87
Case Study: Watermarks in Apache Flink	88
Case Study: Source Watermarks for Google Cloud Pub/Sub	90
Summary	93
<b>4. Advanced Windowing.....</b>	<b>95</b>
<i>When/Where:</i> Processing-Time Windows	95
Event-Time Windowing	97
Processing-Time Windowing via Triggers	98
Processing-Time Windowing via Ingress Time	100
<i>Where:</i> Session Windows	103
<i>Where:</i> Custom Windowing	107
Variations on Fixed Windows	108
Variations on Session Windows	115
One Size Does Not Fit All	119
Summary	119
<b>5. Exactly-Once and Side Effects.....</b>	<b>121</b>
Why Exactly Once Matters	121
Accuracy Versus Completeness	122
Side Effects	123
Problem Definition	123
Ensuring Exactly Once in Shuffle	125
Addressing Determinism	126
Performance	127
Graph Optimization	127
Bloom Filters	128
Garbage Collection	129
Exactly Once in Sources	130
Exactly Once in Sinks	131

Use Cases	133
Example Source: Cloud Pub/Sub	133
Example Sink: Files	134
Example Sink: Google BigQuery	135
Other Systems	136
Apache Spark Streaming	136
Apache Flink	136
Summary	138

---

## Part II. Streams and Tables

<b>6. Streams and Tables.....</b>	<b>141</b>
Stream-and-Table Basics Or: a Special Theory of Stream and Table Relativity	142
Toward a General Theory of Stream and Table Relativity	143
Batch Processing Versus Streams and Tables	144
A Streams and Tables Analysis of MapReduce	144
Reconciling with Batch Processing	150
<i>What, Where, When, and How</i> in a Streams and Tables World	150
<i>What</i> : Transformations	150
<i>Where</i> : Windowing	154
<i>When</i> : Triggers	157
<i>How</i> : Accumulation	165
A Holistic View of Streams and Tables in the Beam Model	166
A General Theory of Stream and Table Relativity	171
Summary	172
<b>7. The Practicalities of Persistent State.....</b>	<b>175</b>
Motivation	175
The Inevitability of Failure	176
Correctness and Efficiency	177
Implicit State	178
Raw Grouping	179
Incremental Combining	181
Generalized State	184
Case Study: Conversion Attribution	186
Conversion Attribution with Apache Beam	189
Summary	199
<b>8. Streaming SQL.....</b>	<b>201</b>
What Is Streaming SQL?	201
Relational Algebra	202

Time-Varying Relations	203
Streams and Tables	207
Looking Backward: Stream and Table Biases	214
The Beam Model: A Stream-Biased Approach	214
The SQL Model: A Table-Biased Approach	218
Looking Forward: Toward Robust Streaming SQL	226
Stream and Table Selection	227
Temporal Operators	228
Summary	249
<b>9. Streaming Joins.....</b>	<b>253</b>
All Your Joins Are Belong to Streaming	253
Unwindowed Joins	254
FULL OUTER	255
LEFT OUTER	258
RIGHT OUTER	259
INNER	259
ANTI	261
SEMI	262
Windowed Joins	266
Fixed Windows	267
Temporal Validity	269
Summary	282
<b>10. The Evolution of Large-Scale Data Processing.....</b>	<b>283</b>
MapReduce	284
Hadoop	288
Flume	289
Storm	294
Spark	297
MillWheel	300
Kafka	304
Cloud Dataflow	307
Flink	309
Beam	313
Summary	316
<b>Index.....</b>	<b>319</b>