

AIQL: Enabling Efficient Attack Investigation from System Monitoring Data

Peng Gao¹, Xusheng Xiao², Zhichun Li³, Kangkook Jee², Fengyuan Xu³, Sanjeev R. Kulkarni¹, Prateek Mittal¹

Advanced Persistent Threat (APT) Attack

APT attacks have plagued many well-protected businesses

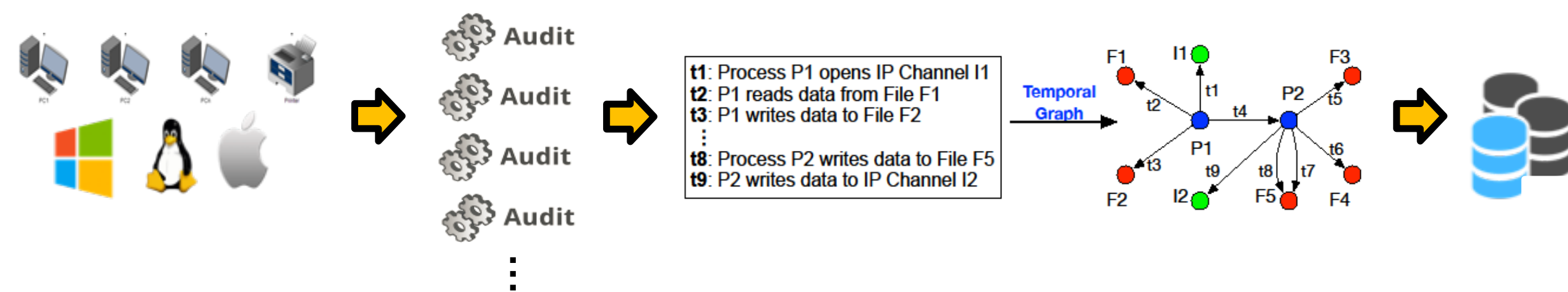


- **Advanced:** sophisticated techniques exploiting multiple vulnerabilities
- **Persistent:** continuously monitoring and stealing data from target
- **Threat:** strong economical or political motives

Ubiquitous System Monitoring

System monitoring records system events from kernels in a unified structure of logs (not bound to applications)

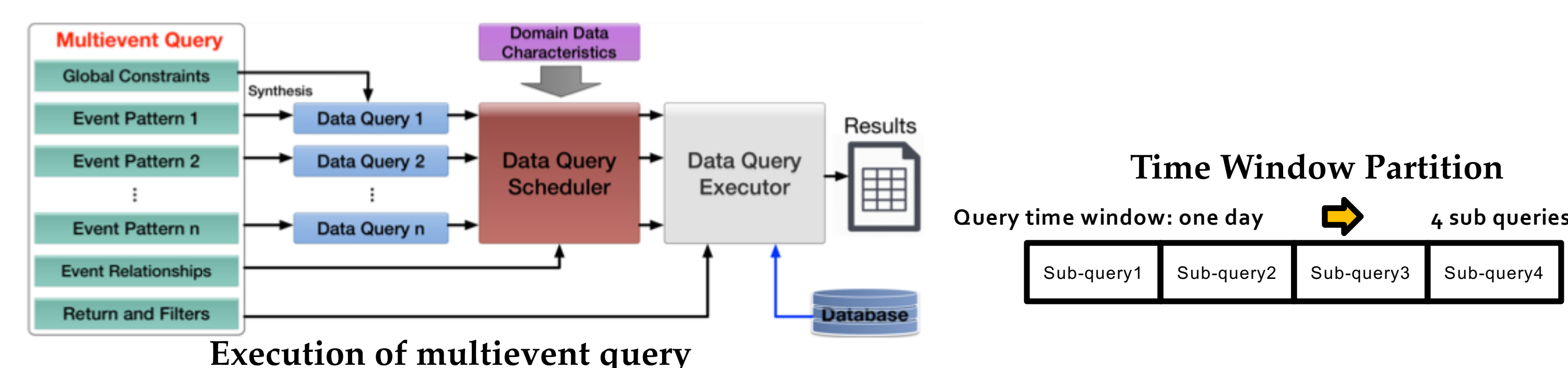
- System activities (system events): <subject, operation, object>



Challenges in enabling efficient attack investigation

- Attack behavior specification
- Timely “big data” security analysis
 - Collect and store system monitoring data for hosts in an organization (~50 GB for 100 hosts per day) => **data storage optimization**
 - Query data for attack investigation => **query execution optimization**

Query Execution Engine

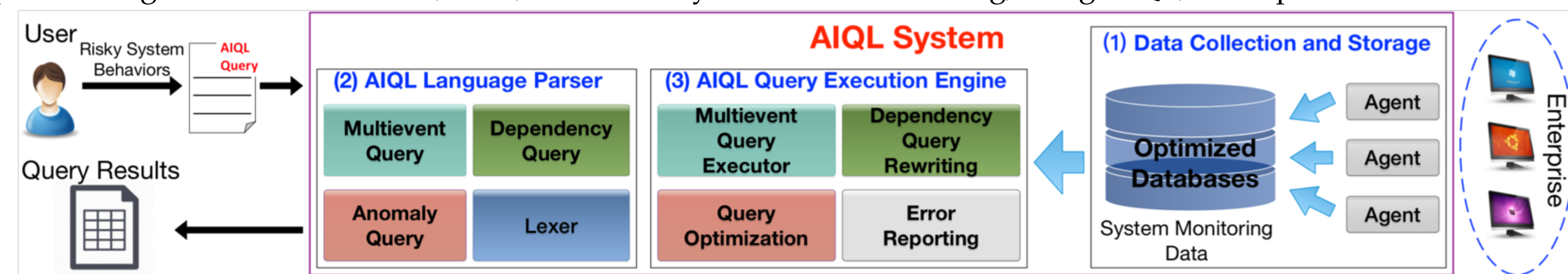


- Synthesize a SQL data query for every event pattern
- Schedule the data queries using domain-specific optimizations
 - Leverage **event relationships** for optimizing **search strategies**
 - **Prioritize** event search based on estimated pruning power
 - **Prune** search space of related events
 - Leverage **domain-specific characteristics** of data for **parallel search**
 - Time window partition

AIQL System

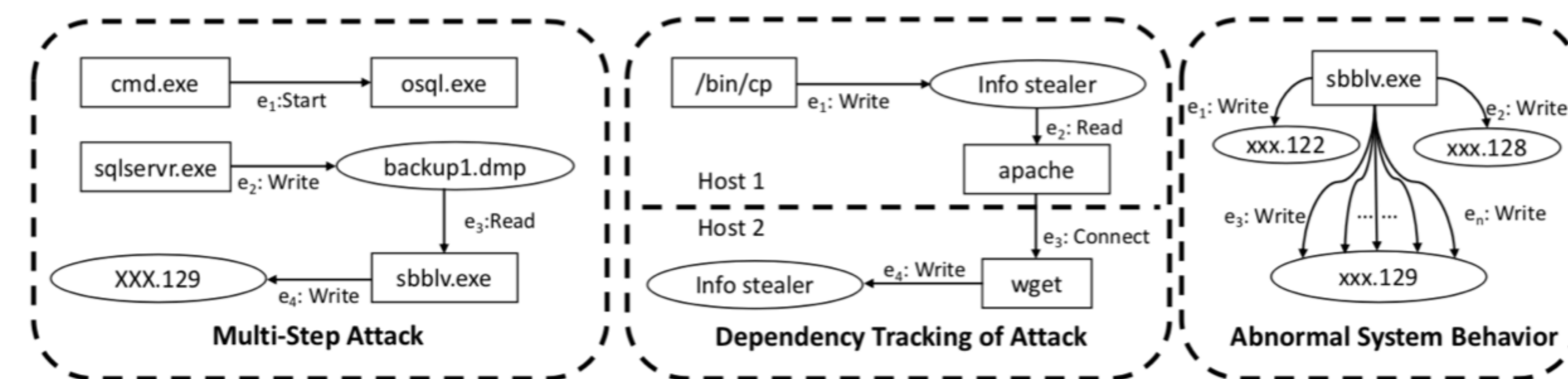
Novel query system for attack investigation (~50K lines of Java code)

- Built on top of existing mature tools: auditd, ETW, DTrace for system-level monitoring; PostgreSQL, Greenplum for relational databases



AIQL (Attack Investigation Query Language) Design

AIQL supports the specification of three types of attack behaviors. Typical constructs: event patterns (subject-operation-object), global constraints, attribute/temporal relationships, constraint chaining, sliding windows and history state access, syntax shortcuts (e.g., attribute inference)



```
1 (at "mm/dd/yyyy")
2 agentid = xxx // SQL database server (obfuscated)
3 proc p1["cmd.exe"] start proc p2["osql.exe"] as
  evt1
4 proc p3["sqlservr.exe"] write file f1["%backp1.dmp"
  ] as evt2
5 proc p4["ssblv.exe"] read file f1 as evt3
6 proc p4 read || write ip i1[dstip="XXX.129"] as evt4
7 with evt1 before evt2, evt2 before evt3, evt3 before
  evt4
8 return distinct p1, p2, p3, p4, i1
```

Multievent AIQL query

```
1 (at "01/01/2017")
2 forward: proc p1["%/bin/cp%", agentid = 2] ->[write]
3     file f1["var/www/info_stealer%"]
4 <-[read] proc p2["%apache%"]
5 ->[connect] proc p3[agentid=3] // tracking across
6     host
7 ->[write] file f2["%info_stealer%"]
8 return f1, p1, p2, p3, f2
```

Dependency AIQL query

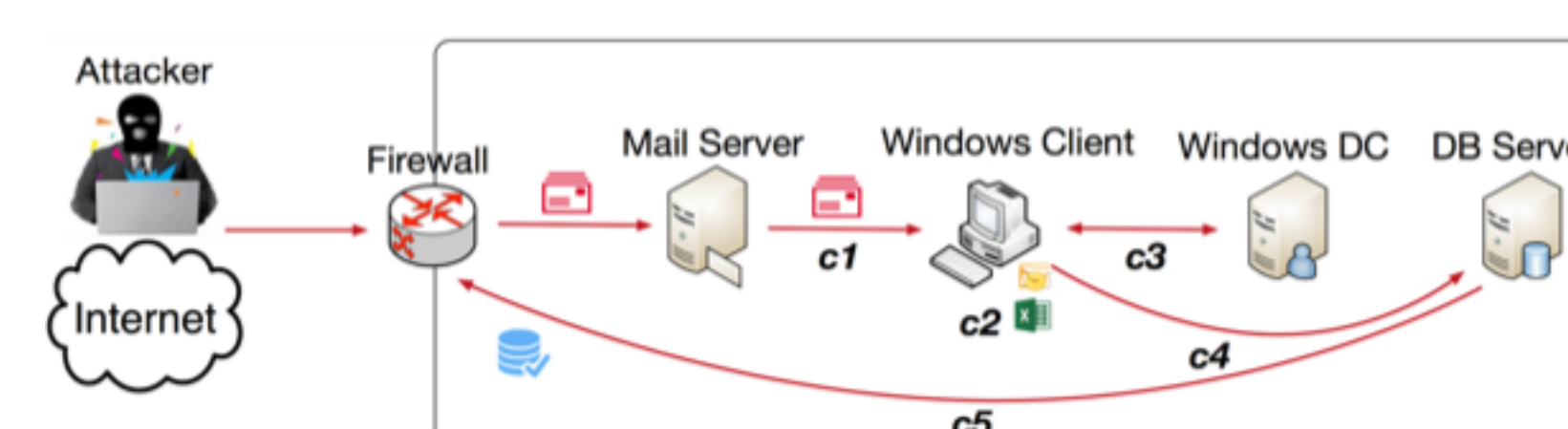
```

1 (at "mm/dd/2017") // date (obfuscated)
2 agentid = xxx // SQL database server (obfuscated)
3 window = 1 min, step = 10 sec
4 proc p write ip i[datip="XXX.129"] as evt
5 return p, avg(evt.amount) as amt
6 group by p
7 having (amt > 2 * (amt + amt[1] + amt[2]) / 3)

```

Anomaly AIQL query

Case Study: APT Attack Investigation



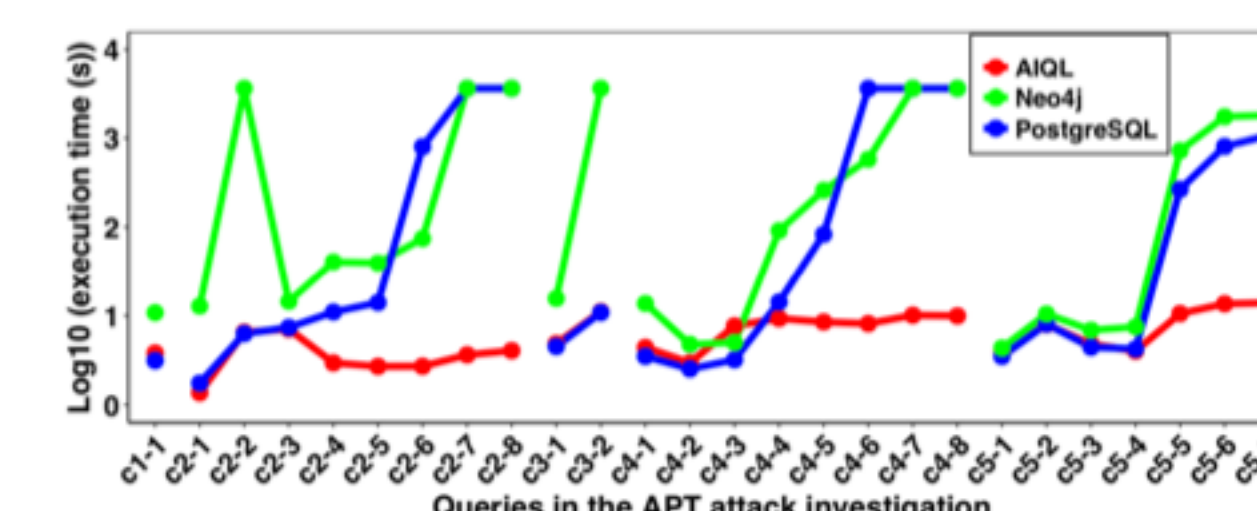
- **c1** Initial Compromise
- **c2** Malware Infection
- **c3** Privilege Escalation
- **c4** Penetration into Database Server
- **c5** Data Exfiltration

27 queries, touching 119 GB of data (422 million system events)

- As the attack behaviors become more complex, SQL and Cypher queries become verbose with many joins and constraints
- AIQL takes less than **3 minutes** for the entire investigation process , **127x** faster than PostgreSQL, **157x** faster than Neo4j

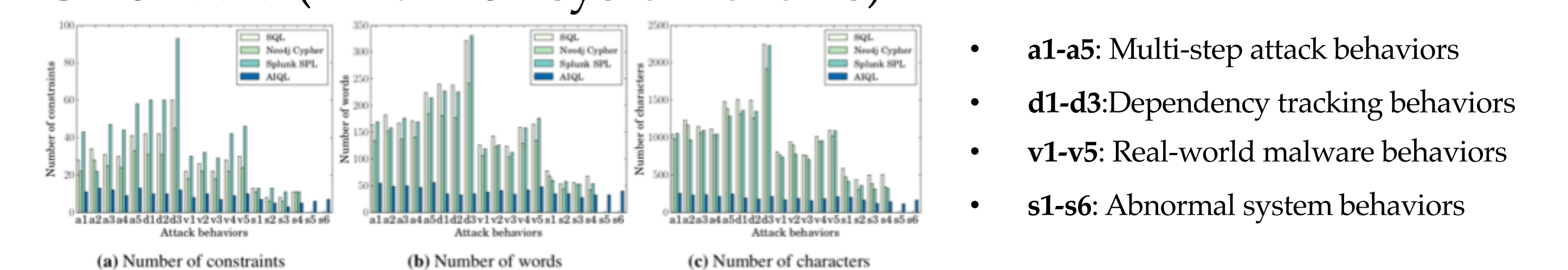
Table 3: Aggregate statistics for case study

Attack type	# of Queries	# of Evt Patterns	AIQL (s)	PostgreSQL (s)	Neo4j (s)
c1	1	3	3.8	3.1	10.8
c2	8	27	31.0	8038.7	10981.7
c3	2	4	15.9	15.3	3615.6
c4	8	35	61.0	10906.7	8150.6
c5	7	18	58.8	2166.5	4285.4
All	26	87	170.5	21130.3	27044.1



Conciseness and Scheduling Efficiency

19 queries on four major types of attack behaviors, touching 738 GB of data (2.1 billion system events)



Other languages vs. AIQL: **2.4x** more constraints, **3.1x** more words, **4.7x** more characters

- Both PostgreSQL and Greenplum employ our **optimized** data storage

