

Contents

1	Appendix C: Production Cost Analysis and Economic Viability	1
1.1	C.1 Cost Modeling Methodology	1
1.1.1	C.1.1 Assumptions	1
1.1.2	C.1.2 Instance Sizing	1
1.2	C.2 Private Information Retrieval (PIR) Costs	2
1.2.1	C.2.1 Computational PIR (Single-Server)	2
1.2.2	C.2.2 Information-Theoretic PIR (Multi-Server)	3
1.3	C.3 Zero-Knowledge Proof Costs	4
1.3.1	C.3.1 Simple Proofs (15K Constraints)	4
1.3.2	C.3.2 Complex Proofs (1M Constraints)	5
1.4	C.4 Combined Stack Costs	6
1.4.1	C.4.1 Small Clinical Practice	6
1.4.2	C.4.2 Research Institution	6
1.4.3	C.4.3 Healthcare Network	7
1.5	C.5 Cost Optimization Strategies	8
1.5.1	C.5.1 Proof Caching	8
1.5.2	C.5.2 Batch Processing	9
1.5.3	C.5.3 Reserved Instances	9
1.6	C.6 Break-Even Analysis	9
1.6.1	C.6.1 CPIR vs IT-PIR	9
1.6.2	C.6.2 Groth16 vs Halo2	10
1.7	C.7 TCO (Total Cost of Ownership) Comparison	10
1.7.1	C.7.1 3-Year TCO Analysis	10
1.7.2	C.7.2 On-Premise vs Cloud	11
1.8	C.8 Pricing Calculator	12
1.9	C.9 Summary & Recommendations	13
1.9.1	C.9.1 Cost Summary	13
1.9.2	C.9.2 Recommendations	13

1 Appendix C: Production Cost Analysis and Economic Viability

1.1 C.1 Cost Modeling Methodology

1.1.1 C.1.1 Assumptions

Pricing Base: AWS us-east-1 (January 2025) - Representative of major cloud providers - On-demand pricing (conservative estimate) - Regional variations: $\pm 15\%$

Query Volume: 10,000 queries/day = 300,000/month - Typical for research institution - Small clinic: ~1,000 queries/day - Large healthcare network: ~100,000 queries/day

Cost Components: 1. **Fixed:** Instance costs (monthly) 2. **Variable:** Per-query compute + network + storage 3. **One-time:** Setup costs (ceremonies, deployment)

1.1.2 C.1.2 Instance Sizing

Compute Requirements:

Component	CPU Intensity	Memory Usage	Instance Type
HDC Encoding	Low (1.49ms)	100MB	t3.small
ZK Proving (15K)	Medium (603ms)	4.2GB	c5.xlarge
ZK Proving (1M)	High (11.2s)	48GB	c5.9xlarge
PIR CPIR (100K)	Medium (590ms)	1.2GB	t3.medium
PIR IT-PIR (3-server)	High (6.4s \times 3)	3.6GB	3 \times m5.xlarge

Burst Credit Analysis (t3 instances):

t3.medium baseline:

Baseline: $20\% \times 2 \text{ vCPU} = 0.4 \text{ vCPU continuous}$

Credits earned: 24 credits/hour

PIR workload:

- 10K queries/day = 417/hour
- CPU time: $417 \times 0.59\text{s} = 246 \text{ CPU-seconds/hour}$
- Credits consumed: $246/3600 \times 60 = 4.1 \text{ credits/hour}$
- Net: +19.9 credits/hour (SUSTAINABLE)

t3.large baseline:

Baseline: $30\% \times 2 \text{ vCPU} = 0.6 \text{ vCPU continuous}$

Credits earned: 36 credits/hour

ZK workload (simple):

- 10K proofs/day = 417/hour
- CPU time: $417 \times 1.15\text{s} = 480 \text{ CPU-seconds/hour}$
- Credits consumed: $480/3600 \times 60 = 8 \text{ credits/hour}$
- Net: +28 credits/hour (SUSTAINABLE)

Conclusion: t3 instances viable for sustained 10K queries/day workload, but c5/m5 recommended for predictable performance.

1.2 C.2 Private Information Retrieval (PIR) Costs

1.2.1 C.2.1 Computational PIR (Single-Server)

100K Database:

Variable per query:

Compute: $0.59\text{s} \times \$0.042/3600 = \0.0000069

Network: $0.0001\text{GB} \times \$0.09 = \0.0000090

Total variable: \$0.0000159

Monthly (300K queries):

Variable: $300,000 \times \$0.000016 = \4.80

Fixed (t3.medium @ \$0.042/hr): \$30

Total: \$35/month

Per Query: \$0.000117

1M Database:

Variable per query:

Compute: $0.92s \times \$0.042/3600 = \0.0000107

Network: $0.001GB \times \$0.09 = \0.0000900

Total variable: $\$0.0001007$

Monthly (300K queries):

Variable: $300,000 \times \$0.00010 = \30

Fixed (t3.large @ $\$0.042/hr$): $\$61$

Total: $\$91/month$

Per Query: $\$0.000303$

10M Database (monolithic):

Variable per query:

Compute: $113s \times \$0.192/3600 = \0.006027

Network: $0.01GB \times \$0.09 = \0.000900

Total variable: $\$0.006927$

Monthly (300K queries):

Variable: $300,000 \times \$0.00693 = \$2,079$

Fixed (r5.xlarge @ $\$0.192/hr$): $\$183$

Total: $\$2,262/month$

Per Query: $\$0.00754$

10M Database (sharded $10 \times 1M$):

Sharding Strategy:

Hash-based routing: $query_hash \% 10 \rightarrow$ single shard

Each shard: 1M records, receives 1K queries/day

Per Shard:

Monthly cost: $\$91$ (same as 1M monolithic at 10K/day)

Total:

$10 \text{ shards} \times \$91 = \$910/month$

Savings vs monolithic: $\$2,262 - \$910 = \$1,352/month$ (60%)

Per Query: $\$0.00303$

1.2.2 C.2.2 Information-Theoretic PIR (Multi-Server)

100K Database (3-server):

Variable per query (3 servers):

Compute: $6.4s \times 3 \times \$0.042/3600 = \0.000224

Network: $0.000538GB \times \$0.09 = \0.000048

Total variable: $\$0.000272$

Monthly (300K queries):

Variable: $300,000 \times \$0.00027 = \81

Fixed (3 \times t3.large @ \$0.042/hr each): \$183

Total: \$264/month

Per Query: \$0.000880

Trust Model: Information-theoretic (requires 2+ honest servers)

1M Database (3-server):

Variable per query (3 servers):

Compute: $8.1s \times 3 \times \$0.096/3600 = \0.000648

Network: $0.0054GB \times \$0.09 = \0.000486

Total variable: \$0.001134

Monthly (300K queries):

Variable: $300,000 \times \$0.00113 = \339

Fixed (3 \times m5.xlarge @ \$0.096/hr each): \$415

Total: \$754/month

Per Query: \$0.002513

Trust Model: Information-theoretic (unconditional privacy)

1.3 C.3 Zero-Knowledge Proof Costs

1.3.1 C.3.1 Simple Proofs (15K Constraints)

Groth16:

Proving:

Time: 1.15s per proof

Instance: c5.large (2 vCPU, 4GB)

Hourly rate: \$0.085

Proofs per hour: 3,130 (theoretical), 939 (30% util)

Variable per proof:

Compute: $1.15s \times \$0.085/3600 = \0.000027

Network: $192B \times \$0.09/1GB = \0.000000

Storage: negligible

Total: \$0.000027

Monthly (300K proofs):

Variable: $300,000 \times \$0.000027 = \8.10

Fixed (c5.large): \$61

Total: \$69/month

Setup (one-time): \$10-50K (ceremony)

PLONK:

Proving:

Time: 0.82s per proof
Instance: c5.xlarge (4 vCPU, 8GB)
Hourly rate: \$0.170
Proofs per hour: 4,390 (theoretical), 1,317 (30% util)

Variable per proof:

Compute: $0.82s \times \$0.170/3600 = \0.000031
Network: $1KB \times \$0.09/1GB = \0.000000
Total: \$0.000031

Monthly (300K proofs):

Variable: $300,000 \times \$0.000031 = \9.30
Fixed (c5.xlarge): \$122
Total: \$131/month

Setup (one-time): \$0 (use universal SRS)

Halo2 (Recommended):

Proving:

Time: 0.60s per proof
Instance: c5.xlarge (4 vCPU, 8GB)
Hourly rate: \$0.170
Proofs per hour: 6,000 (theoretical), 1,800 (30% util)

Variable per proof:

Compute: $0.60s \times \$0.170/3600 = \0.000028
Network: $5KB \times \$0.09/1GB = \0.000000
Total: \$0.000028

Monthly (300K proofs):

Variable: $300,000 \times \$0.000028 = \8.40
Fixed (c5.xlarge): \$122
Total: \$130/month

Setup (one-time): \$0 (trustless)

1.3.2 C.3.2 Complex Proofs (1M Constraints)

Halo2 (Recommended for complex):

Proving:

Time: 11.2s per proof
Instance: c5.9xlarge (36 vCPU, 72GB)
Hourly rate: \$1.530
Proofs per hour: 321 (theoretical), 96 (30% util)

Variable per proof:

Compute: $11.2s \times \$1.530/3600 = \0.004760

Network: $5KB \times \$0.09/1GB = \0.000000

Total: $\$0.004760$

Monthly (300K proofs):

Variable: $300,000 \times \$0.00476 = \$1,428$

Fixed (c5.9xlarge): $\$1,101$

Total: $\$2,529/\text{month}$

Peak memory: 48GB

1.4 C.4 Combined Stack Costs

1.4.1 C.4.1 Small Clinical Practice

Configuration:

Scale: 1,000 patients

Query volume: 10,000 queries/day

Components:

- PIR: CPIR, 100K database
- ZK: Halo2, 15K constraints (simple queries)
- HDC: On-demand encoding

Instances:

- t3.medium (PIR): $\$30/\text{month}$
- c5.xlarge (ZK): $\$122/\text{month}$
- t3.small (HDC/API): $\$15/\text{month}$

Performance:

- PIR latency: 590ms
- ZK latency: 600ms
- Total E2E: $\sim 1.2s$

Monthly Cost:

PIR: $\$35$

ZK: $\$132$

Total: $\$167/\text{month}$

Per Query: $\$0.000556$

Comparison:

Traditional cloud genomics: $\$3,000\text{--}5,000/\text{month}$

Savings: 95%

1.4.2 C.4.2 Research Institution

Configuration:

Scale: 100,000 samples

Query volume: 10,000 queries/day

Components:

- PIR: IT-PIR (3-server), 1M database
- ZK: Halo2, 15K constraints
- HDC: Batch encoding

Instances:

- 3 × m5.xlarge (PIR): \$415/month
- c5.xlarge (ZK): \$122/month
- t3.large (HDC/API): \$61/month

Performance:

- PIR latency: 8.1s (IT-PIR)
- ZK latency: 600ms
- Total E2E: ~8.7s

Trust Model:

- PIR: Information-theoretic (2+ honest servers)
- ZK: Trustless (Halo2, no ceremony)

Monthly Cost:

PIR: \$754
ZK: \$132
Total: \$886/month

Per Query: \$0.00295

Comparison:

Traditional platform: \$5,000–8,000/month
Savings: 85%

1.4.3 C.4.3 Healthcare Network

Configuration:

Scale: 10M records

Query volume: 10,000 queries/day

Components:

- PIR: CPIR sharded (10 × 1M), hash routing
- ZK: Halo2, 1M constraints (complex PRS)
- HDC: Distributed encoding

Instances:

- 10 × t3.large (PIR shards): \$910/month
- c5.9xlarge (ZK): \$1,101/month
- c5.2xlarge (HDC/API): \$549/month

Performance:

- PIR latency: 920ms (single shard)
- ZK latency: 11.2s (complex proof)
- Total E2E: ~12.1s

Sharding Strategy:

Hash-based: `query_hash % 10`
Load balanced: ~1K queries/shard/day
Fault tolerance: Each shard can handle 10K/day

Monthly Cost:

PIR (sharded): \$910
ZK (complex): \$2,529
Total: \$3,439/month

Per Query: \$0.01146

Comparison:

Traditional platform: \$15,000-30,000/month
Savings: 77%

Note: Sharding reduces PIR cost by 60% vs monolithic

1.5 C.5 Cost Optimization Strategies

1.5.1 C.5.1 Proof Caching

Implementation:

Cache Layer: Redis cluster

Strategy:

- Cache proven queries
- TTL: 24 hours
- Max size: 10GB
- Eviction: LRU

Hit Rates (measured):

- Variant presence: 42%
- PRS queries: 18%
- Ancestry checks: 65%
- Overall: 40%

Cost Impact:

Before caching: \$132/month (ZK)
After caching (40% hits): \$79/month
Savings: \$53/month (40%)

Redis cost: \$15/month (elasticache.t3.small)

Net savings: \$38/month (29%)

1.5.2 C.5.2 Batch Processing

Implementation:

Strategy:

- Queue queries during day
- Batch process at night (off-peak)
- Use spot instances (70% discount)

Instance Selection:

On-demand c5.9xlarge: \$1.530/hr
Spot c5.9xlarge: \$0.459/hr (70% discount)

Batch Efficiency:

Serial: 321 proofs/hour (100% util)
Parallel: 2,568 proofs/hour (8 instances)

Cost Comparison (300K proofs/month):

On-demand (30% util): \$2,529/month
Spot batch (80% util): \$892/month
Savings: \$1,637/month (65%)

Trade-off: Higher latency (overnight processing)

1.5.3 C.5.3 Reserved Instances

3-Year Reserved Savings:

Instance Type: c5.xlarge

On-demand: \$0.170/hr = \$122/month

1-year reserved: \$0.111/hr = \$80/month (35% savings)

3-year reserved: \$0.084/hr = \$60/month (51% savings)

Upfront Payment (3-year):

No upfront: \$60/month × 36 = \$2,160 total

All upfront: \$1,825 total (15% additional savings)

Recommended: 3-year all-upfront for stable workloads

1.6 C.6 Break-Even Analysis

1.6.1 C.6.1 CPIR vs IT-PIR

Formula:

$$Q^* = (F - F) / (30 \cdot (v - v))$$

where: - F = fixed monthly cost - v = variable cost per query - Q^* = break-even queries/day

100K Database:

CPIR: F = \$30, v = \$0.000016

IT-PIR: $F = \$183$, $v = \$0.00027$

$Q^* = (183 - 30) / (30 \times (0.000016 - 0.00027))$

$Q^* = 153 / -0.00762$

$Q^* = -20,079$ queries/day

Since $Q^* < 0$, CPIR is ALWAYS cheaper

(Higher fixed costs of IT-PIR never recover from lower variable costs)

Conclusion: Choose IT-PIR for **unconditional privacy**, not cost. CPIR is always more economical but requires computational assumptions.

1.6.2 C.6.2 Groth16 vs Halo2

15K Constraints:

Groth16: $F = \$61$, $v = \$0.000027$

Halo2: $F = \$122$, $v = \$0.000028$

$Q^* = (122 - 61) / (30 \times (0.000027 - 0.000028))$

$Q^* = 61 / -0.00003$

$Q^* = -2,033,333$ queries/day

Since $Q^* < 0$, Halo2 is NEVER cheaper operationally

However: Halo2 avoids \$10-50K trusted setup ceremony. Break-even:

Setup cost savings: \$30K (typical)

Monthly premium: $\$122 - \$61 = \$61$

Months to amortize: $\$30,000 / \$61 = 492$ months (41 years)

Conclusion: Choose Halo2 for **trustless** security, not cost. Setup avoidance worth the premium for regulatory/audit reasons.

1.7 C.7 TCO (Total Cost of Ownership) Comparison

1.7.1 C.7.1 3-Year TCO Analysis

GenomeVault (Research Institution):

Year 1:

Infrastructure: $\$886/\text{month} \times 12 = \$10,632$

Development: \$50,000 (setup, integration)

Training: \$10,000

Total Year 1: \$70,632

Year 2-3:

Infrastructure: \$10,632/year

Maintenance: \$5,000/year

Total Year 2-3: \$15,632/year

3-Year TCO: \$101,896
Average annual: \$33,965

Traditional Cloud Genomics Platform:

Year 1:

Platform fees: \$6,000/month × 12 = \$72,000
Setup/migration: \$25,000
Training: \$15,000
Total Year 1: \$112,000

Year 2-3:

Platform fees: \$72,000/year
Support: \$10,000/year
Total Year 2-3: \$82,000/year

3-Year TCO: \$276,000
Average annual: \$92,000

Savings: \$174,104 over 3 years (63% reduction)

1.7.2 C.7.2 On-Premise vs Cloud

On-Premise Hardware:

Initial Investment:

Servers (3× Dell R750): \$45,000
Storage (NAS 100TB): \$30,000
Networking: \$15,000
Setup/installation: \$20,000
Total: \$110,000

Annual Operating:

Power (3kW @ \$0.12/kWh): \$3,154
Cooling (additional 1kW): \$1,051
Maintenance: \$15,000
Staff (0.5 FTE): \$50,000
Total annual: \$69,205

3-Year TCO: \$110,000 + 3×\$69,205 = \$317,615
Average annual: \$105,872

Cloud (GenomeVault):

3-Year TCO: \$101,896
Average annual: \$33,965

Conclusion: Cloud is 68% cheaper than on-premise for typical research institution scale. On-premise becomes competitive at >100K queries/day scale.

1.8 C.8 Pricing Calculator

Interactive Calculator (Python):

```
def calculate_monthly_cost(
    queries_per_day: int,
    database_rows: int,
    zk_constraints: int = 15000,
    backend: str = "halo2",
    pir_type: str = "cpir"
) -> dict:
    # PIR costs
    pir_config = {
        ("cpir", 100_000): (30, 0.000016),
        ("cpir", 1_000_000): (61, 0.00010),
        ("cpir", 10_000_000): (183, 0.00693),
        ("itpir", 100_000): (183, 0.00027),
        ("itpir", 1_000_000): (415, 0.00113),
    }

    # ZK costs
    zk_config = {
        ("halo2", 15_000): (122, 0.000028),
        ("plonk", 15_000): (122, 0.000031),
        ("groth16", 15_000): (61, 0.000027),
        ("halo2", 1_000_000): (1101, 0.004760),
    }

    # Get costs
    pir_fixed, pir_var = pir_config.get((pir_type, database_rows), (61, 0.0001))
    zk_fixed, zk_var = zk_config.get((backend, zk_constraints), (122, 0.000028))

    # Calculate monthly
    queries_per_month = queries_per_day * 30

    pir_variable_monthly = pir_var * queries_per_month
    zk_variable_monthly = zk_var * queries_per_month

    total_monthly = (
        pir_fixed + pir_variable_monthly +
        zk_fixed + zk_variable_monthly
    )

    cost_per_query = (pir_variable_monthly + zk_variable_monthly) / queries_per_month

    return {
        "monthly_cost": total_monthly,
        "cost_per_query": cost_per_query,
```

```

        "breakdown": {
            "pir_fixed": pir_fixed,
            "pir_variable": pir_variable_monthly,
            "zk_fixed": zk_fixed,
            "zk_variable": zk_variable_monthly,
        }
    }

# Example usage
cost = calculate_monthly_cost(
    queries_per_day=10_000,
    database_rows=1_000_000,
    zk_constraints=15_000,
    backend="halo2",
    pir_type="cpir"
)

print(f"Monthly cost: ${cost['monthly_cost']:.2f}")
print(f"Per query: ${cost['cost_per_query']:.6f}")

```

1.9 C.9 Summary & Recommendations

1.9.1 C.9.1 Cost Summary

Deployment	Monthly Cost	Per Query	vs Traditional	Best For
Small Clinic	\$167	\$0.00056	95% savings	1K patients
Research Institution	\$886	\$0.00295	85% savings	100K samples
Healthcare Network	\$3,439	\$0.01146	77% savings	10M records

1.9.2 C.9.2 Recommendations

Start with: - CPIR for performance - Halo2 for trustless proofs - t3/c5 instances for predictable costs

Optimize with: - Proof caching (40% cost reduction) - Batch processing (65% reduction for flexible workloads) - Reserved instances (50% reduction for stable loads)

Scale with: - PIR sharding (60% savings at 10M+ scale) - Multi-region deployment (20% premium for redundancy) - Spot instances (70% discount for batch workloads)

Calculator Tool: Interactive cost calculator available at `scripts/cost_calculator.py`