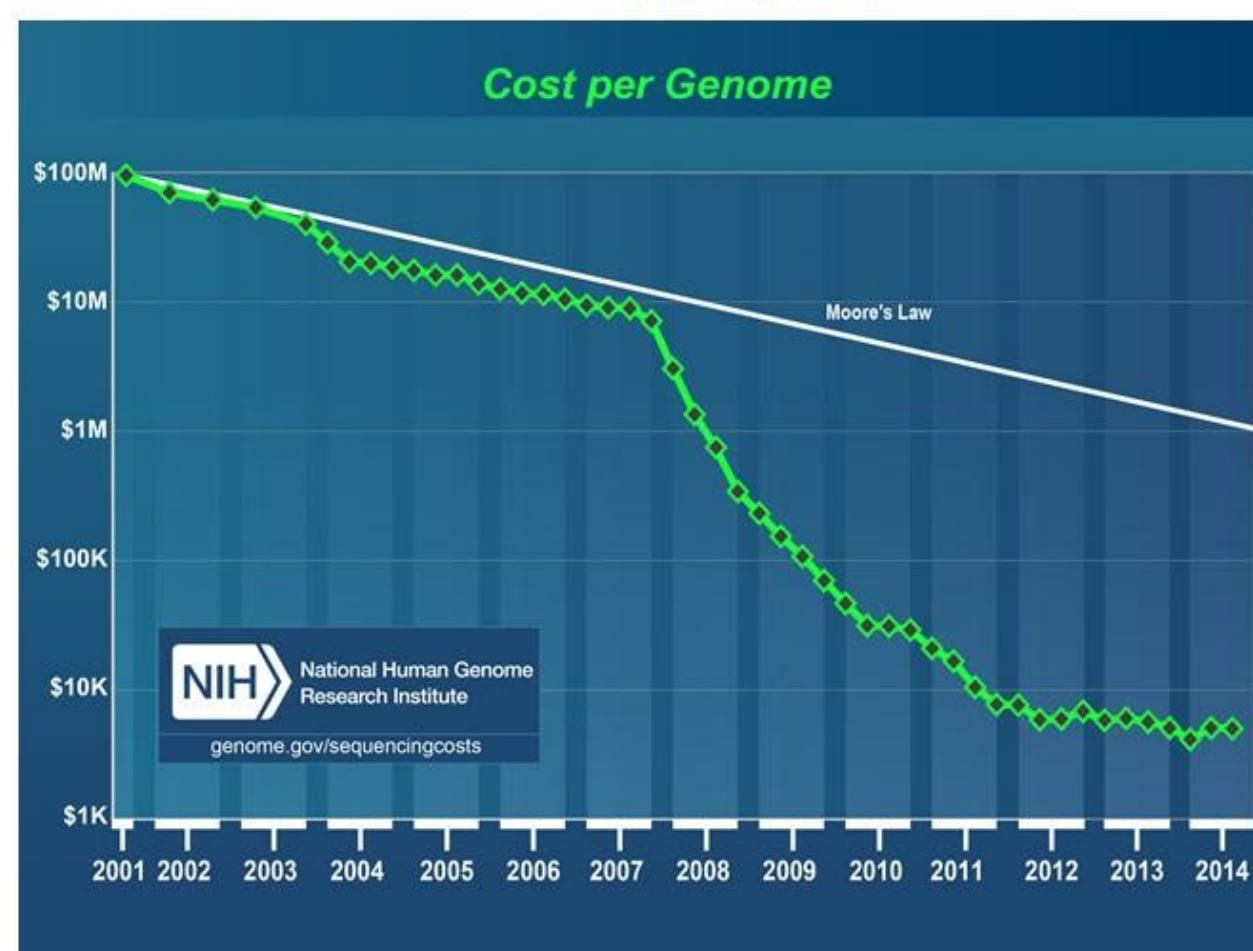


Context



Vraag naar efficiënte, schaalbare analyse en opslag van gegevens ↗



Oplossing

• Dataschema

- Genotype-kolommen => kolom voor elke sample

variant_id	ref	alt	...	gt_type_alex	gt_type_john	...	gt_depth_alex	gt_depth_john	...
------------	-----	-----	-----	--------------	--------------	-----	---------------	---------------	-----

- Arbitraire queries:

- Zonder indices
- Duplicatie data: extra tabellen met geschikte keys definiëren

sample_name	gt_depth	variant_id	chrom	start	variant_id	sex	sample_name
-------------	----------	------------	-------	-------	------------	-----	-------------

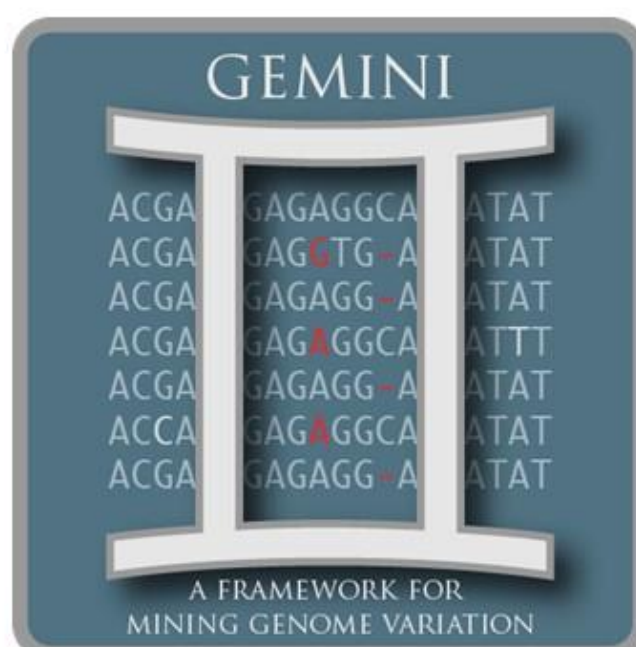
- Extra tabellen verbergen voor gebruiker
- Geen joins => denormalisatie

• Inladen genomdata

- Vullen variants-tabel efficiënt te paralleliseren
- Cores schrijven parallel nr 1 Cassandra cluster
- >> SQLite: elk nr 1 DB, achteraf mergen



Context: case study



- Multi-sample genomanalyse
- Annotaties, vb:
 - aaf_1kg_afr (frequentie genetische variant in Afrikaanse populatie)
 - is_somatic (somatische mutatie, mogelijk carcinogeen)
- Ad-hoc queries
- SQLite = te traag, schaal niet

```
$ gemini query --header \
-q "SELECT chrom, start, end, ref, alt, gene, (gts).(*) \
FROM variants" \
--gt-filter "(gt_types).(phenotype==1).(==HOM_REF).(all) \
and \
(gt_depths).(phenotype==2).(==HOM_REF).(none)" \
extended_ped.db
```



Oplossing

• Queries

- Query op hoofdtabel (variants, samples) => splitsen in subqueries op hulptabellen. Bvb:

```
SELECT * FROM variants WHERE chrom = 'chromX'
AND start > 1500 AND gt_type_alex = HET
=> SELECT variant_id FROM variants_by_chrom_start
WHERE chrom = 'chromX' AND start > 1500
+ SELECT variant_id FROM variants_by_samples_gt_type
WHERE sample_name = 'Alex' AND gt_type = HET
```

- Resultaten subqueries = verzamelingen

- Combineren met set-operaties

- Resulterende rijen, kolommen

- opvragen uit hoofdtabel

Query	Resultaat
$p \text{ AND } r$	$res(p) \cap res(r)$
$p \text{ OR } r$	$res(p) \cup res(r)$
$\text{NOT } p$	$I \setminus res(p)$



Probleem

- Relatieel dataschema GEMINI

- variants-tabel
 - Genotype-kolommen: compressed binary vectors
- samples-table: geslacht, herkomst, familierelaties

variant_id	chrom	start	...	gts	...
1	chromX	24	...	A/A A/G T/T T/C
2	chromY	3541	...	T/T T/T T/C A/G

- Ad-hoc queries:

- Standaard SQL-queries (onmogelijk op binary genotype-arrays)
- Daarom: --gt-filter: gt_type.alex == HET
- Wildcards: (COL).(SAMPLE_WILDCARD).(WILDCARD_RULE).(RULE_ENFORCEMENT)
Bvb. (gt_types).(sex==1).(==HET).(all)



Oplossing

• --gt-filter wildcards

- Lange ketens subqueries op genotype-hulptabellen
- Evaluatie subqueries voor verschillende samples => paralleliseren

• User interface

- Gebruiker moet hulptabellen kennen
- SQL-syntaxis licht gewijzigd => parsen
 - Queries binnen 1 hulptabel: CQL syntaxis
 - Gekoppeld met && i.p.v. AND, ||, NOT

```
$ gemini query -q "SELECT chrom, start, subtype FROM variants
WHERE chrom = 'chromX' AND start > 5600 AND gene = 'gene A'
--gt-filter "(gt_types).(phenotype==1).(==HOM_REF).(all)"
```

```
$ geminicassandra query -q "SELECT chrom, start, subtype FROM variants
WHERE chrom = 'chromX' AND start > 5600 && gene = 'gene A'
--gt-filter "(gt_types).(phenotype==1).(==HOM_REF).(all)"
```



Apache Cassandra i.p.v. SQLite

- Schaalbaarheid ✓
- Relatiele tabellen => columnair model ✓
- Concurrency ✓
- CQL ~ SQL; Python API ✓

Maar...

- Minder support voor indexing, querying
 - Geen joins
 - Beperkingen op WHERE-clause
 - » Cassandra = gigantische multimap <keys, columns> -> values
 - » Niet itereren en alle rijen filteren, maar hash(key) berekenen -> map nr correcte rijen
 - » Enkel opeenvolgend opgeslagen rijen opvragen
 - » PRIMARY KEY ((col1,...), (col_x)*)
partition key | clustering key
bepaalt node in cluster | bepaalt volgorde rijen met == partition key
 - » Geen OR, NOT, !=, ...



Conclusies

- Doel: schaalbare genomanalyse
- GEMINI = genoom-DB + uitgebreide ad-hoc SQL-queries
- Apache Cassandra: compromis
 - Schaalbaarheid vs. query-features
 - Dataduplicatie, denormalisatie
 - Queries in applicatie-laag
 - Parallelisatie
 - Flexibel schema

