

RELATIONAL DATA MODELS FOR GENETIC VCF DATA

Mohamed Sabri Hafidi, Ozan Kahramanoğulları, Anton Dignös, and Johann Gamper

{mhmohamedsabri, okahramanogullari, anton.dignoes, johann.gamper}@unibz.it

CHALLENGES

- Genomic data is growing exponentially, reaching petabyte volumes.
- The structure of VCF is compact but not optimized for advanced querying.
- Existing tools struggle with large-scale datasets and complex analyses.
- Merging VCF data with diverse biological sources remains limited.

OBJECTIVES

- Population-level storage
- Flexible SQL-based querying
- Optimized performance
- Seamless multi-omics integration

RELATED WORK

- Command-line tools:** Offer speed but lack flexibility. e.g., *BCFtools*, *Tabix*.
- Libraries:** Enable custom analyses but require multiple tools for a complete workflow. e.g., *vcflib*, *cvcf2*.
- Database systems:** Scale well but vary in how they represent the VCF data. e.g., *TileDB-VCF*, *TheSNPpit*.
- Prior studies explore *relational* and *JSON-based* models but lack broad evaluations.
- We introduce and evaluate novel relational models for VCF data management.

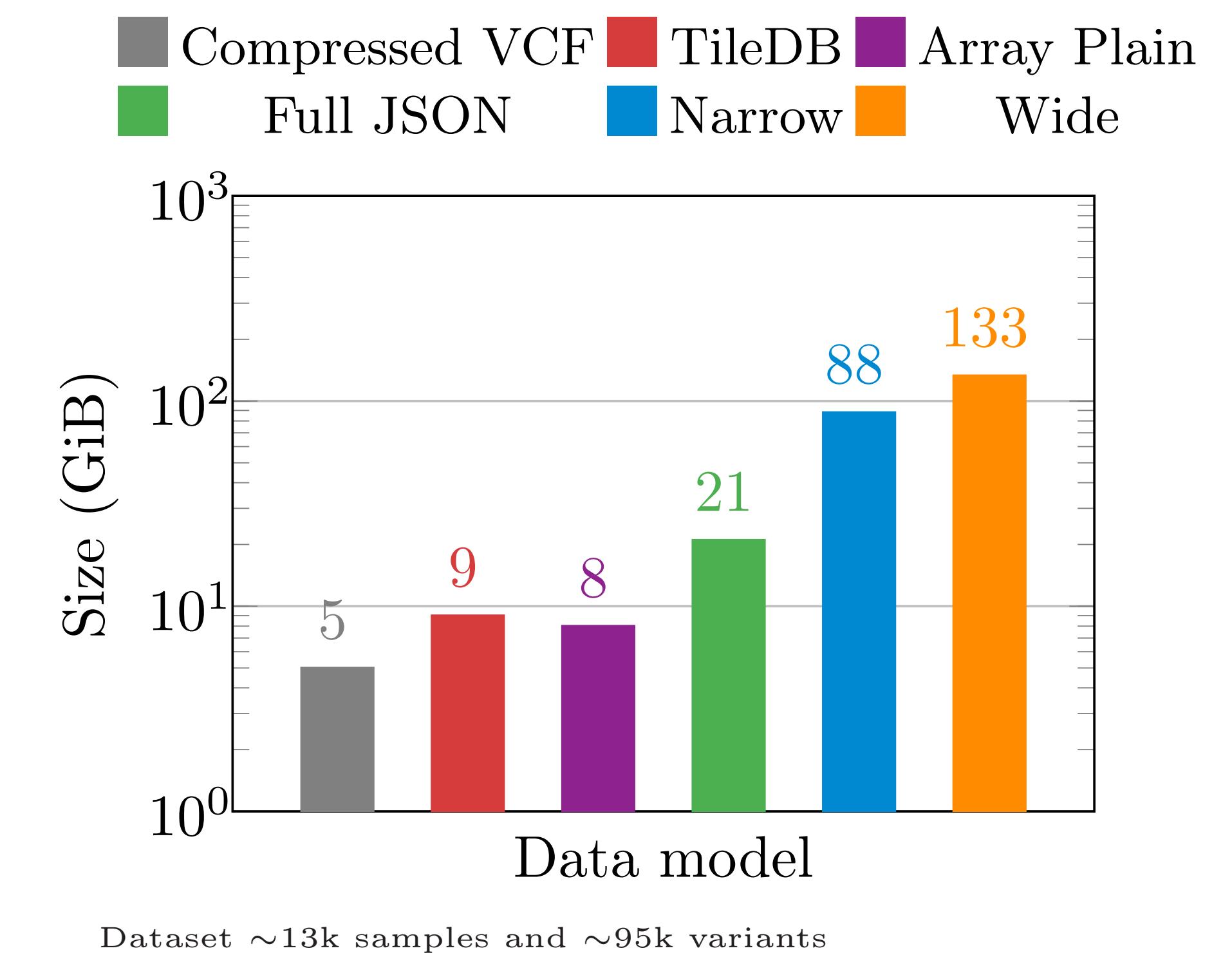
THE VCF FILE STRUCTURE

| | | Meta-Information | | | | | | | | | |
|----------|--------------|------------------|---------|---------|----|------|-----------------------------------|-------------|----------------|----------------|--------------|
| | | Header | | | | | | | | | |
| Variants | Fixed Fields | Samples | | | | | | | | | |
| | | NA00001 | NA00002 | NA00003 | | | | | | | |
| 20 | 14370 | rs6054257 | G | A | 29 | PASS | NS=3;DP=14;AF=0.5;DB:H2 | GT:GQ:DP:HQ | 0 0:48:1:51,51 | 1 0:48:8:51,51 | 1/1:43:5:... |
| 20 | 17330 | . | T | A | 3 | q10 | NS=3;DP=11;AF=0.017 | GT:GQ:DP:HQ | 0 0:49:3:58,50 | 0 1:3:5:65,3 | 0/0:41:3 |
| 20 | 1110696 | rs6040355 | A | G,T | 67 | PASS | NS=2;DP=10;AF=0.333,0.667;AA=T;DB | GT:GQ:DP:HQ | 1 2:21:6:23,27 | 2 1:2:0:18,2 | 2/2:35:4 |
| 20 | 1230237 | . | T | . | 47 | PASS | NS=3;DP=13;AA=T | GT:GQ:DP:HQ | 0 0:54:7:56,60 | 0 0:48:4:51,51 | 0/0:61:2 |
| 20 | 1234567 | microsat1 | GTC | G,GTCT | 50 | PASS | NS=3;DP=9;AA=G | GT:GQ:DP | 0 1:35:4 | 0/2:17:2 | 1/1:40:3 |

RELATIONAL DATA MODELS

| Wide: | each row represents a variant and each column represents a sample | | | | | | | | | | | | | | | | | | |
|--------------------------|---|---|-----------|-----------|---|-----------|---|---------|---|---|---|-----------|---|---|-----------|---|---|-----------|--|
| vcf_wide_samples_chunk_1 | <table border="1"> <tr><th>ln</th><th>s_id</th><th>genotypes</th></tr> <tr><td>1</td><td>NA00001</td><td>0 0:48:1:51,51</td></tr> <tr><td>2</td><td>NA00002</td><td>1 0:48:8:51,51</td></tr> <tr><td>3</td><td>NA00003</td><td>1/1:43:5</td></tr> <tr><td>4</td><td>NA00004</td><td>2 0:49:3:58,50</td></tr> <tr><td>5</td><td>NA00005</td><td>0 0:54:7:56,60</td></tr> </table> | ln | s_id | genotypes | 1 | NA00001 | 0 0:48:1:51,51 | 2 | NA00002 | 1 0:48:8:51,51 | 3 | NA00003 | 1/1:43:5 | 4 | NA00004 | 2 0:49:3:58,50 | 5 | NA00005 | 0 0:54:7:56,60 |
| ln | s_id | genotypes | | | | | | | | | | | | | | | | | |
| 1 | NA00001 | 0 0:48:1:51,51 | | | | | | | | | | | | | | | | | |
| 2 | NA00002 | 1 0:48:8:51,51 | | | | | | | | | | | | | | | | | |
| 3 | NA00003 | 1/1:43:5 | | | | | | | | | | | | | | | | | |
| 4 | NA00004 | 2 0:49:3:58,50 | | | | | | | | | | | | | | | | | |
| 5 | NA00005 | 0 0:54:7:56,60 | | | | | | | | | | | | | | | | | |
| Narrow: | each row represents a variant-sample pair | | | | | | | | | | | | | | | | | | |
| vcf_narrow_samples | <table border="1"> <tr><th>ln</th><th>s_id</th><th>genotypes</th></tr> <tr><td>1</td><td>NA00001</td><td>0 0:48:1:51,51</td></tr> <tr><td>2</td><td>NA00002</td><td>1 0:48:8:51,51</td></tr> <tr><td>3</td><td>NA00003</td><td>1/1:43:5</td></tr> <tr><td>4</td><td>NA00004</td><td>2 0:49:3:58,50</td></tr> <tr><td>5</td><td>NA00005</td><td>0 0:54:7:56,60</td></tr> </table> | ln | s_id | genotypes | 1 | NA00001 | 0 0:48:1:51,51 | 2 | NA00002 | 1 0:48:8:51,51 | 3 | NA00003 | 1/1:43:5 | 4 | NA00004 | 2 0:49:3:58,50 | 5 | NA00005 | 0 0:54:7:56,60 |
| ln | s_id | genotypes | | | | | | | | | | | | | | | | | |
| 1 | NA00001 | 0 0:48:1:51,51 | | | | | | | | | | | | | | | | | |
| 2 | NA00002 | 1 0:48:8:51,51 | | | | | | | | | | | | | | | | | |
| 3 | NA00003 | 1/1:43:5 | | | | | | | | | | | | | | | | | |
| 4 | NA00004 | 2 0:49:3:58,50 | | | | | | | | | | | | | | | | | |
| 5 | NA00005 | 0 0:54:7:56,60 | | | | | | | | | | | | | | | | | |
| Array Plain: | store the genotypes in an array | | | | | | | | | | | | | | | | | | |
| vcf_array_indices | <table border="1"> <tr><th>s_id</th><th>e_id</th></tr> <tr><td>NA00001</td><td>1</td></tr> <tr><td>NA00002</td><td>2</td></tr> <tr><td>NA00003</td><td>3</td></tr> <tr><td>...</td><td>...</td></tr> </table> | s_id | e_id | NA00001 | 1 | NA00002 | 2 | NA00003 | 3 | ... | ... | | | | | | | | |
| s_id | e_id | | | | | | | | | | | | | | | | | | |
| NA00001 | 1 | | | | | | | | | | | | | | | | | | |
| NA00002 | 2 | | | | | | | | | | | | | | | | | | |
| NA00003 | 3 | | | | | | | | | | | | | | | | | | |
| ... | ... | | | | | | | | | | | | | | | | | | |
| vcf_array_genotypes | <table border="1"> <tr><th>ln</th><th>genotypes</th></tr> <tr><td>1</td><td>"0 0:48:1:51,51", "1 0:48:8:51,51", "1/1:43:5", ...</td></tr> <tr><td>2</td><td>"0 0:49:3:58,50", "1 0:54:7:56,60", "0 0:41:3", ...</td></tr> <tr><td>3</td><td>"1 2:21:6:23,27", "2 1:2:0:18,2", "2 2:35:4", ...</td></tr> <tr><td>4</td><td>"0 0:54:7:56,60", "0 0:48:4:51,51", "0 0:61:2", ...</td></tr> <tr><td>5</td><td>"0 0:135:4", "0/2:17:2", "1/1:40:3", ...</td></tr> </table> | ln | genotypes | 1 | "0 0:48:1:51,51", "1 0:48:8:51,51", "1/1:43:5", ... | 2 | "0 0:49:3:58,50", "1 0:54:7:56,60", "0 0:41:3", ... | 3 | "1 2:21:6:23,27", "2 1:2:0:18,2", "2 2:35:4", ... | 4 | "0 0:54:7:56,60", "0 0:48:4:51,51", "0 0:61:2", ... | 5 | "0 0:135:4", "0/2:17:2", "1/1:40:3", ... | | | | | | |
| ln | genotypes | | | | | | | | | | | | | | | | | | |
| 1 | "0 0:48:1:51,51", "1 0:48:8:51,51", "1/1:43:5", ... | | | | | | | | | | | | | | | | | | |
| 2 | "0 0:49:3:58,50", "1 0:54:7:56,60", "0 0:41:3", ... | | | | | | | | | | | | | | | | | | |
| 3 | "1 2:21:6:23,27", "2 1:2:0:18,2", "2 2:35:4", ... | | | | | | | | | | | | | | | | | | |
| 4 | "0 0:54:7:56,60", "0 0:48:4:51,51", "0 0:61:2", ... | | | | | | | | | | | | | | | | | | |
| 5 | "0 0:135:4", "0/2:17:2", "1/1:40:3", ... | | | | | | | | | | | | | | | | | | |
| Full JSON: | store the sample ID and its genotype in a JSON | | | | | | | | | | | | | | | | | | |
| vcf_json_samples | <table border="1"> <tr><th>ln</th><th>sample</th><th>genotype</th></tr> <tr><td>1</td><td>"NA00001"</td><td>{"id": "NA00001", "genotype": {"DP": "1", "GT": "0 0", "GQ": "48", "HQ": "51,51"}, ...}</td></tr> <tr><td>2</td><td>"NA00002"</td><td>{"id": "NA00002", "genotype": {"DP": "1", "GT": "1 0", "GQ": "49", "HQ": "58,50"}, ...}</td></tr> <tr><td>3</td><td>"NA00003"</td><td>{"id": "NA00003", "genotype": {"DP": "1", "GT": "2 0", "GQ": "54", "HQ": "23,27"}, ...}</td></tr> <tr><td>4</td><td>"NA00004"</td><td>{"id": "NA00004", "genotype": {"DP": "1", "GT": "2 1", "GQ": "52", "HQ": "56,60"}, ...}</td></tr> <tr><td>5</td><td>"NA00005"</td><td>{"id": "NA00005", "genotype": {"DP": "1", "GT": "0 1", "GQ": "35"}, ...}</td></tr> </table> | ln | sample | genotype | 1 | "NA00001" | {"id": "NA00001", "genotype": {"DP": "1", "GT": "0 0", "GQ": "48", "HQ": "51,51"}, ...} | 2 | "NA00002" | {"id": "NA00002", "genotype": {"DP": "1", "GT": "1 0", "GQ": "49", "HQ": "58,50"}, ...} | 3 | "NA00003" | {"id": "NA00003", "genotype": {"DP": "1", "GT": "2 0", "GQ": "54", "HQ": "23,27"}, ...} | 4 | "NA00004" | {"id": "NA00004", "genotype": {"DP": "1", "GT": "2 1", "GQ": "52", "HQ": "56,60"}, ...} | 5 | "NA00005" | {"id": "NA00005", "genotype": {"DP": "1", "GT": "0 1", "GQ": "35"}, ...} |
| ln | sample | genotype | | | | | | | | | | | | | | | | | |
| 1 | "NA00001" | {"id": "NA00001", "genotype": {"DP": "1", "GT": "0 0", "GQ": "48", "HQ": "51,51"}, ...} | | | | | | | | | | | | | | | | | |
| 2 | "NA00002" | {"id": "NA00002", "genotype": {"DP": "1", "GT": "1 0", "GQ": "49", "HQ": "58,50"}, ...} | | | | | | | | | | | | | | | | | |
| 3 | "NA00003" | {"id": "NA00003", "genotype": {"DP": "1", "GT": "2 0", "GQ": "54", "HQ": "23,27"}, ...} | | | | | | | | | | | | | | | | | |
| 4 | "NA00004" | {"id": "NA00004", "genotype": {"DP": "1", "GT": "2 1", "GQ": "52", "HQ": "56,60"}, ...} | | | | | | | | | | | | | | | | | |
| 5 | "NA00005" | {"id": "NA00005", "genotype": {"DP": "1", "GT": "0 1", "GQ": "35"}, ...} | | | | | | | | | | | | | | | | | |

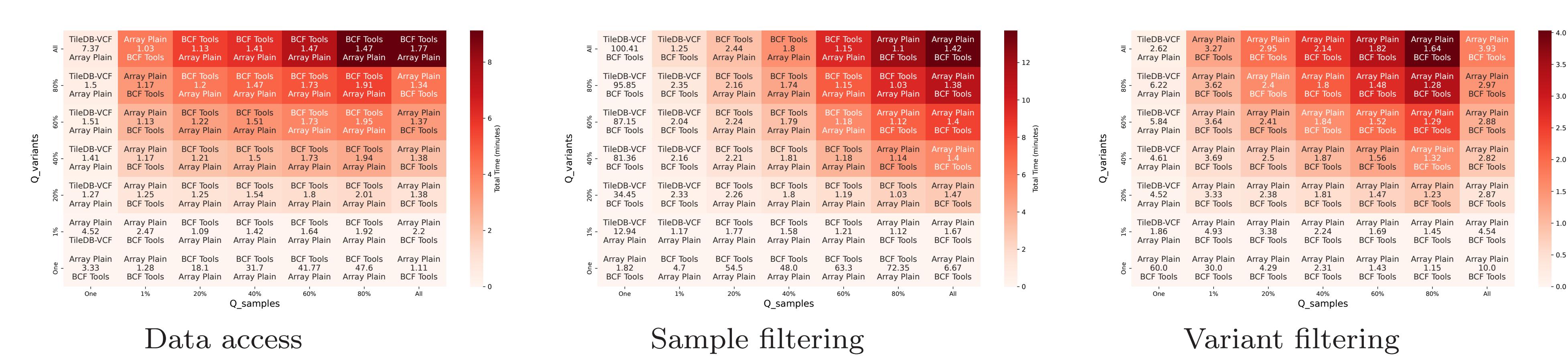
STORAGE REQUIREMENTS



EMPIRICAL EXPERIMENTS

- Data access:** Retrieves entries for a set of sample IDs and/or variant range. “%” denotes a selected amount, and “All” means no filtering is applied.
- Variant filtering:** Extends data access pattern and adds a filter on variant’s fixed fields.
- Sample filtering:** Extends data access pattern and adds a filter on genotype values. Returns matching sample IDs.

Runtime heatmaps, each cell shows the best- and second-best-performing approaches along with their speed-up factor.



SUMMARY

- RDBMS models rival specialized tools.
- No one-size-fits-all solution: Optimal model depends on query workload.
- Relational models offer superior scalability for large VCF datasets.
- Trade-offs are workload-dependent:
 - Storage needs
 - Query selectivity
 - Dataset size
- Array Plain model – Good all-rounder.

Artifact available:

This work was funded in part by the Autonomous Province of Bozen-Bolzano, call Joint Projects South Tyrol – Germany, project DyHealthNet – CUP: I53C22002980003.