

GenomAlx

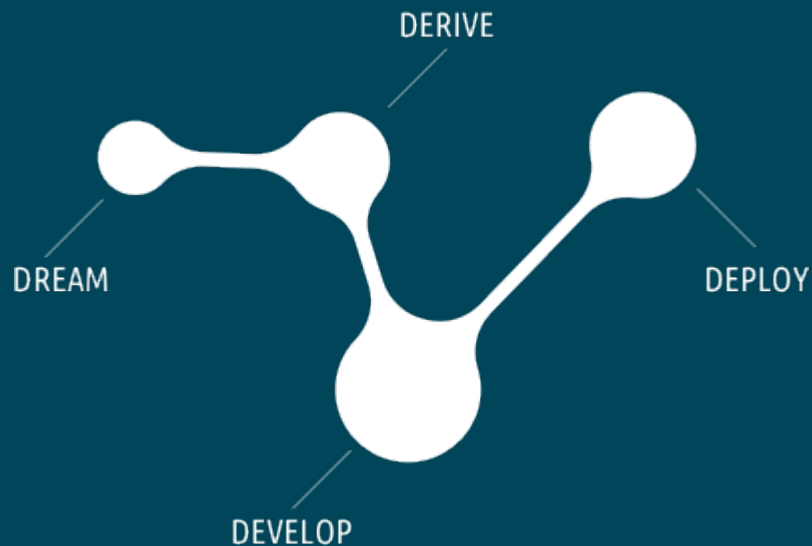
Erforschung des Dark Genome mit KI zur Entwicklung neuartiger
Krankheitsinterventionen bei seltenen Erkrankungen

Outline

1. Short Introduction
2. Motivation
3. Workflow & Cloud Architecture
4. What we achieved so far
5. Next Steps
6. Discussion



Kurzvorstellung: scieneers GmbH



DREAM

Wir starten bei der Beratung, greifen Ihre Ideen auf, entwickeln gemeinsam den Business Case,

DERIVE

entwickeln Datenstrukturen und Modelle,

DEVELOP

implementieren die Lösung

DEPLOY

und integrieren diese nahtlos in Ihre operativen IT-Systeme und Geschäftsprozesse

Wir gewinnen Erkenntnisse aus **Daten** und schaffen damit **Werte**.
Für unsere Kunden, die Gesellschaft und uns selbst.

UNIKLINIK RWTHAACHEN

Institut für Humangenetik
und Genommedizin



GenomAlx

Motivation

Genomsequenzierung bei seltenen Erkrankungen

Seltene Erkrankungen

SELTENE ERKRANKUNGEN: HERAUSFORDERUNG DIAGNOSE

350 Mio., fast 5%
der Weltbevölkerung
leben mit
einer **seltene**
Erkrankung.³



Etwa 75%
der seltenen
Erkrankungen
betreffen
Kinder.¹⁰



Über 7.000 seltene
Erkrankungen
sind **bekannt**.²

Seltene
Erkrankungen
wiegen
schwer

4 Mio.
Betroffene in
Deutschland.⁶



Langer
Weg zur
Diagnose

Typischerweise
werden bis
zu **8 Ärzten**
aufgesucht.³



40% der Patienten
erhalten mindestens
einmal eine
Fehldiagnose.⁷



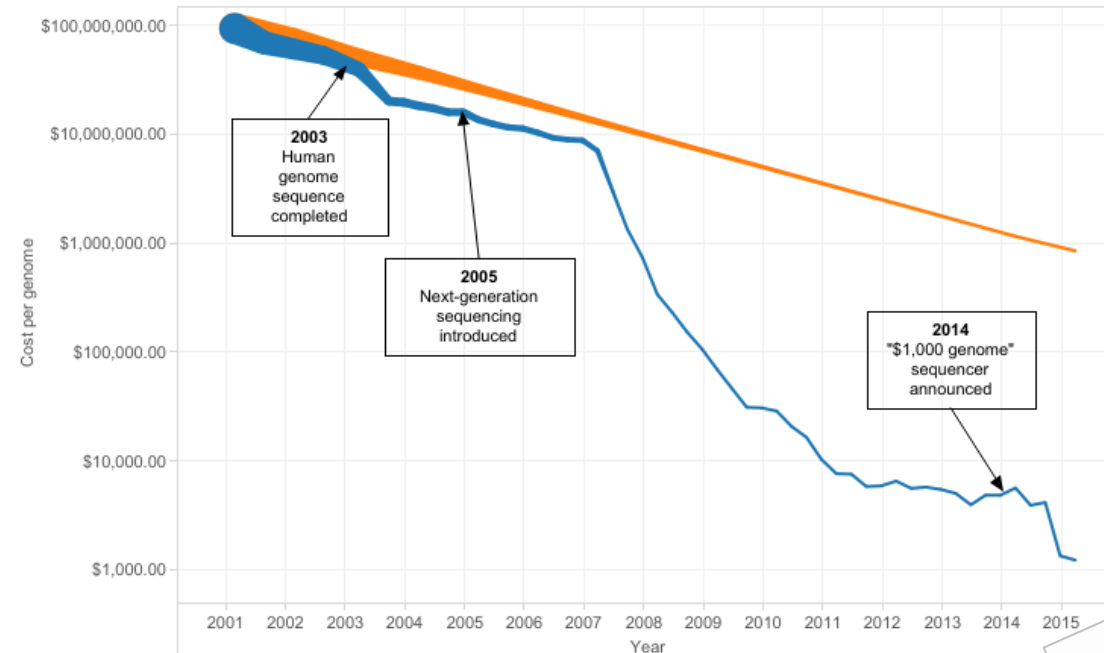
Im Durchschnitt
vergehen bis zur
richtigen Diagnose
4,8 Jahre.³

80% dieser Erkrankungen entstehen durch eine einzige genetische Veränderung und können durch eine Genomanalyse diagnostiziert werden.



Kosten für Genome sinken... Auswertung ist das Bottleneck

DNA sequencing costs over time



Avg. Cost per Mb
\$0,01
\$2.000,00

\$4.000,00
\$6.000,00

Data source
NHGRI data
Moore's law calculation

Decline in real costs compared to expected declines based on Moore's Law.
Trend line: Cost per human genome. Line width: Cost per megabase (Mb)
(Data: NHGRI <https://www.genome.gov/27541954/dna-sequencing-costs-data/>)



Mardis Genome Medicine 2010, 2:84
<http://genomemedicine.com/content/2/11/84>

MUSINGS

The \$1,000 genome, the \$100,000 analysis?

Elaine R Mardis*

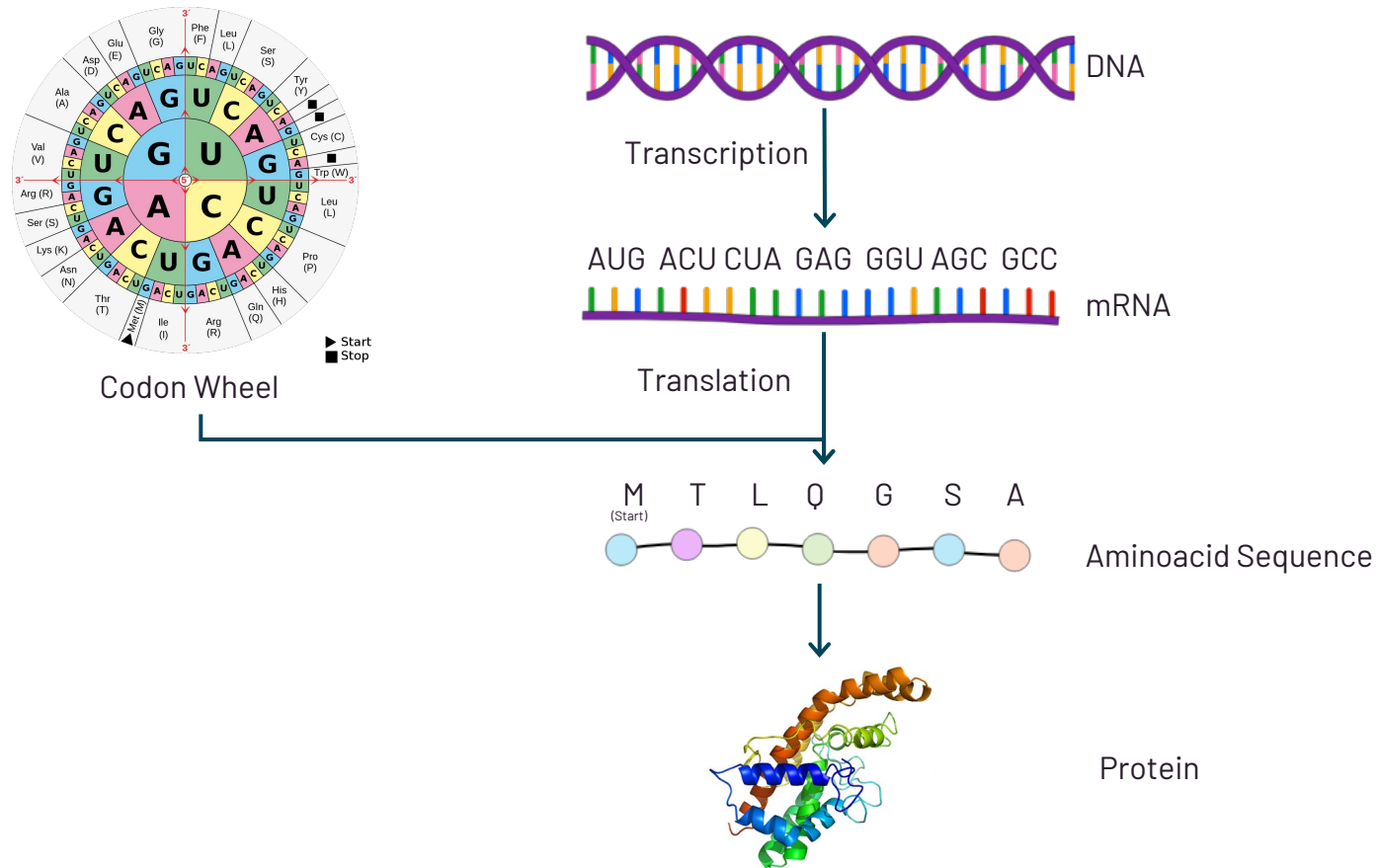
Genomsequenzierung bei seltenen Erkrankungen



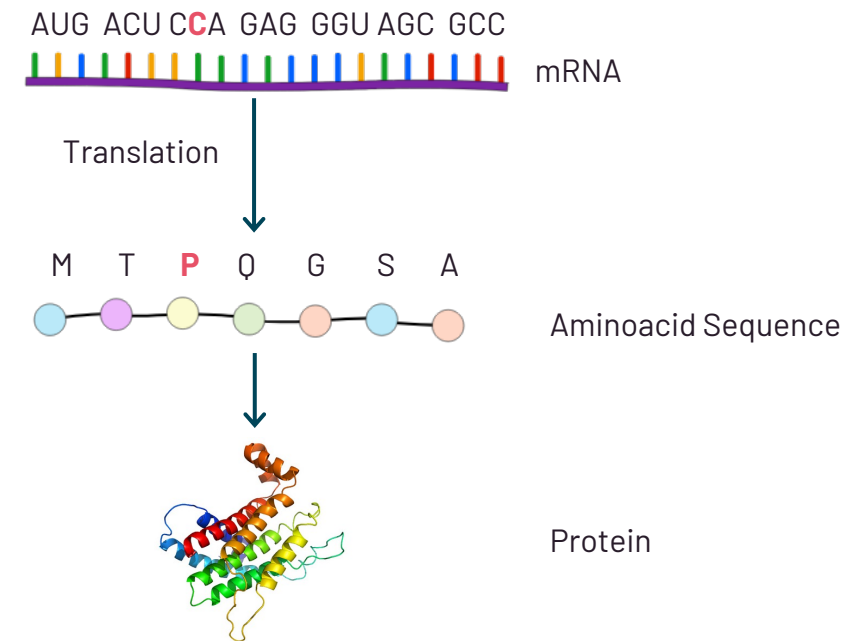
Bausteine des Genoms	3.3 Milliarden
Varianten im Genom	3.5 Millionen
Varianten im Exom	70.000
Seltene Varianten	3.000
„Ultraseltene“ Varianten	300

...but what are variants?

Let's do a short recap on Protein-Synthesis



... but what if we encounter a so called single nucleotide variant (SNV)



Note: On the right side a single nucleotide variant is displayed leading to a **Missense Mutation**. It is important to mention that there are also other types of mutations like **Non-Sense, Synonym and Loss-of-Function Mutations**.

Genomsequenzierung bei seltenen Erkrankungen



Bausteine des Genoms	3.3 Milliarden
Varianten im Genom	3.5 Millionen
Varianten im Exom	70.000
Seltene Varianten	3.000
„Ultraseltene“ Varianten	300

Auf der Suche nach **DER EINEN (!)** Variante, die die Krankheit ausgelöst hat!

Welche Veränderung ist pathogen? Welche harmlos?



Oder...



Verständnis genomischer Varianten

Wo stehen wir?

Gedankenexperiment

Wir bilden die gesamte genetische Information eines Menschen auf 52 Karteikarten ab:

Die Informationen auf **51 Karteikarten entziehen sich** bislang der Bewertbarkeit!

Das entspricht einem **Anteil von 98%** an unserem **Genom**. Dieser Teil des Genoms wird deshalb auch als **Dark Genome** oder **Non-Coding Region** bezeichnet.

Dort sind u.a. die bis zu **300.000** vorhergesagten **small Open Reading Frames** (sORFs) beheimatet. Diese sind bisher kaum verstanden. Einzelne **Varianten in diesen sORFs** konnte bereits eine **Pathogenität nachgewiesen** werden.



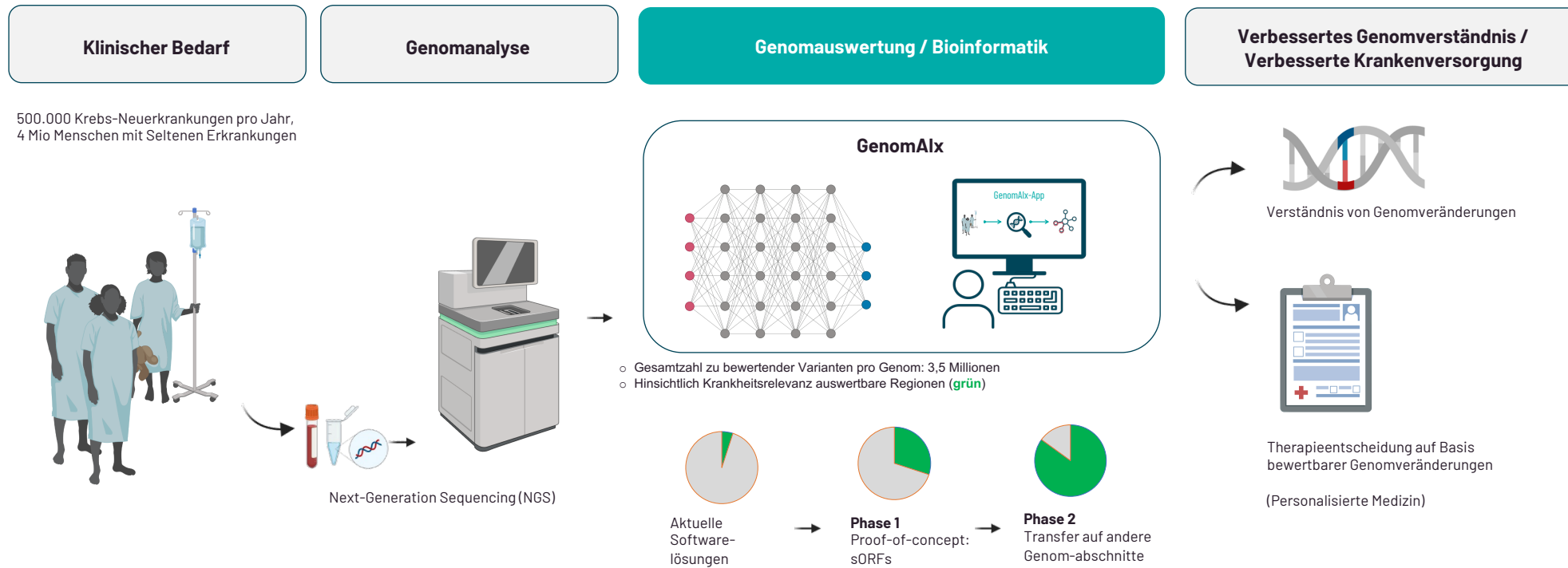
Unser **bisheriges Verständnis** des Genoms kann auf **einer Karteikarte** abgebildet werden.

Das entspricht einem **Anteil von 2%** unseres **Genoms** und wird als **Exom** oder **Coding-Region** bezeichnet

Dort sind u.a. unsere ca. **22.000 Gene** beheimatet, deren Varianten vergleichsweise **gut verstanden** sind.

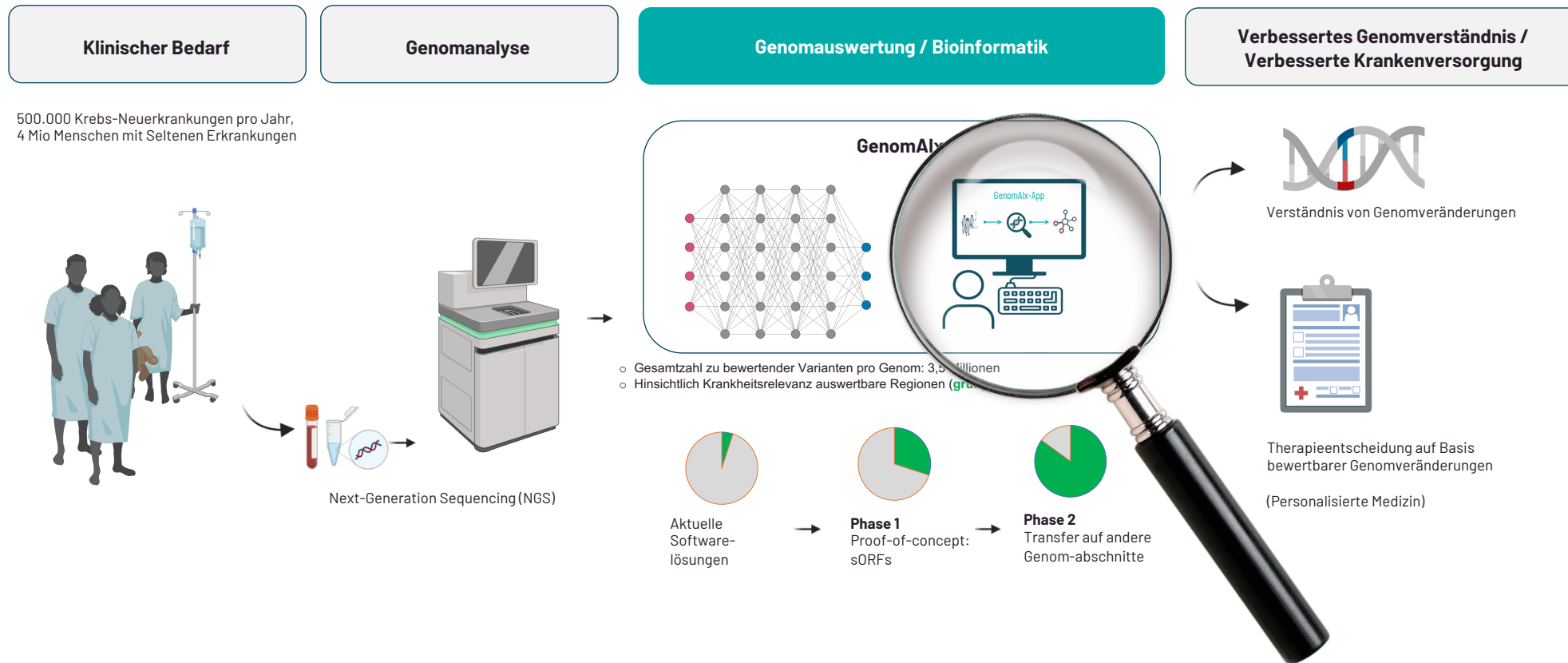
GenomAlx

Erforschung des Dark Genome mit KI zur Entwicklung neuartiger Krankheitsinterventionen



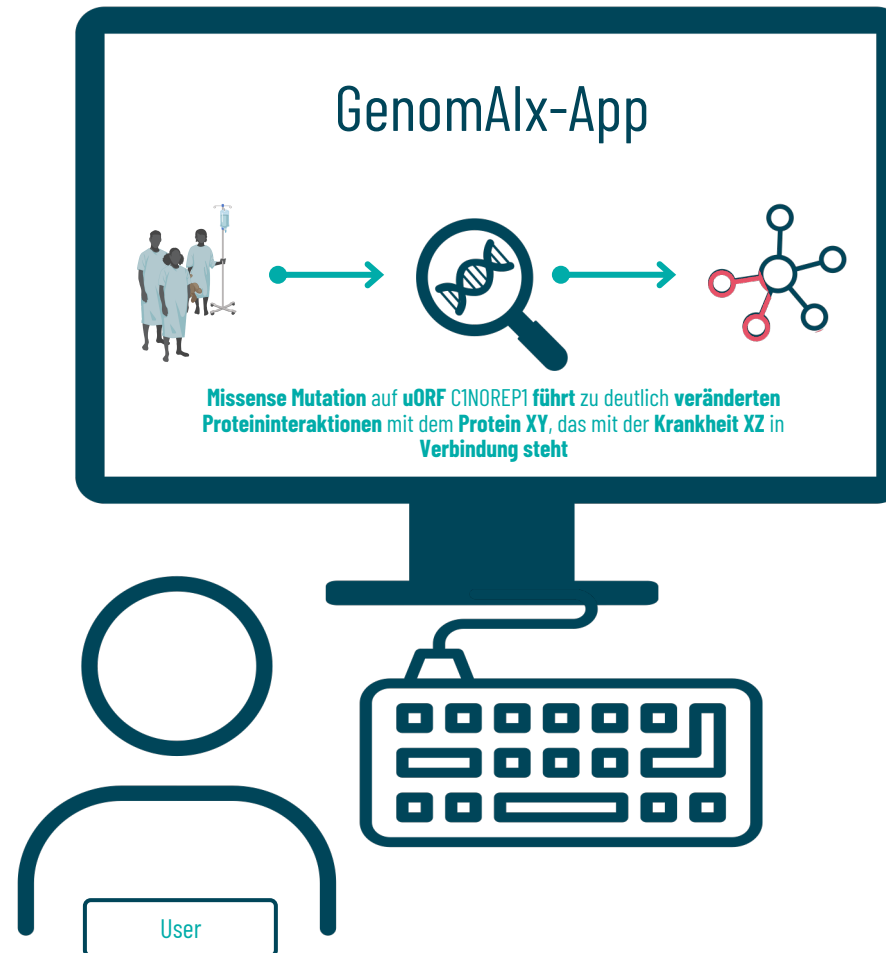
GenomAlx

Erforschung des Dark Genome mit KI zur Entwicklung neuartiger Krankheitsinterventionen



GenomAlx

eine cloudbasierte SaaS-Anwendung zur umfänglicheren Genomauswertung

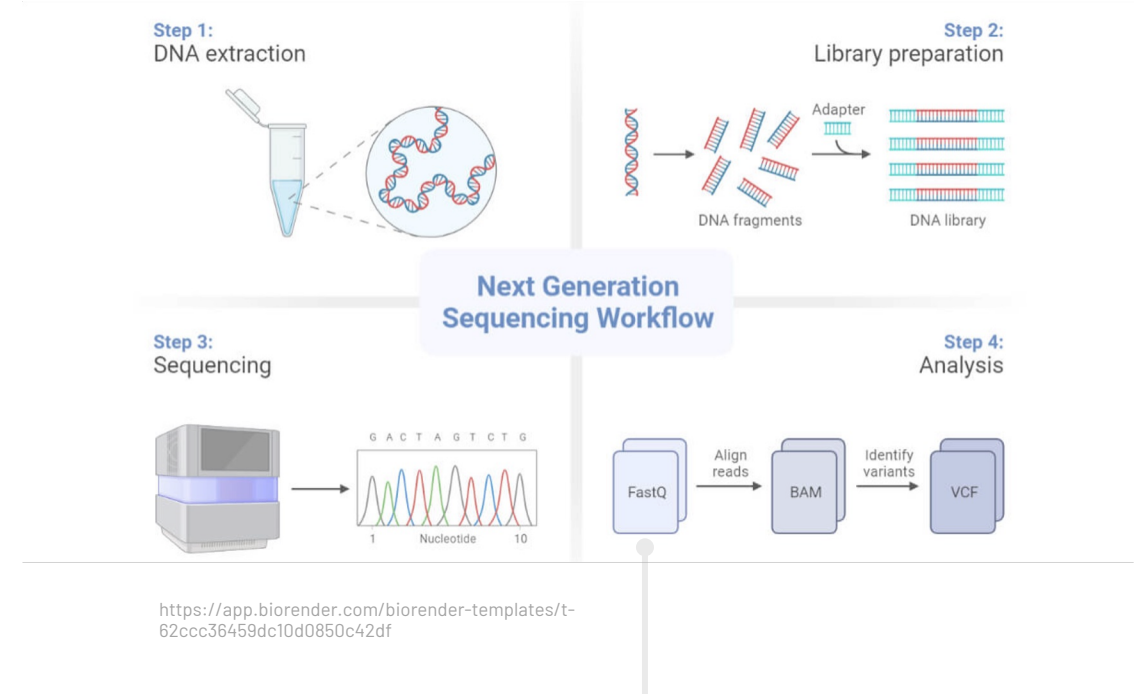
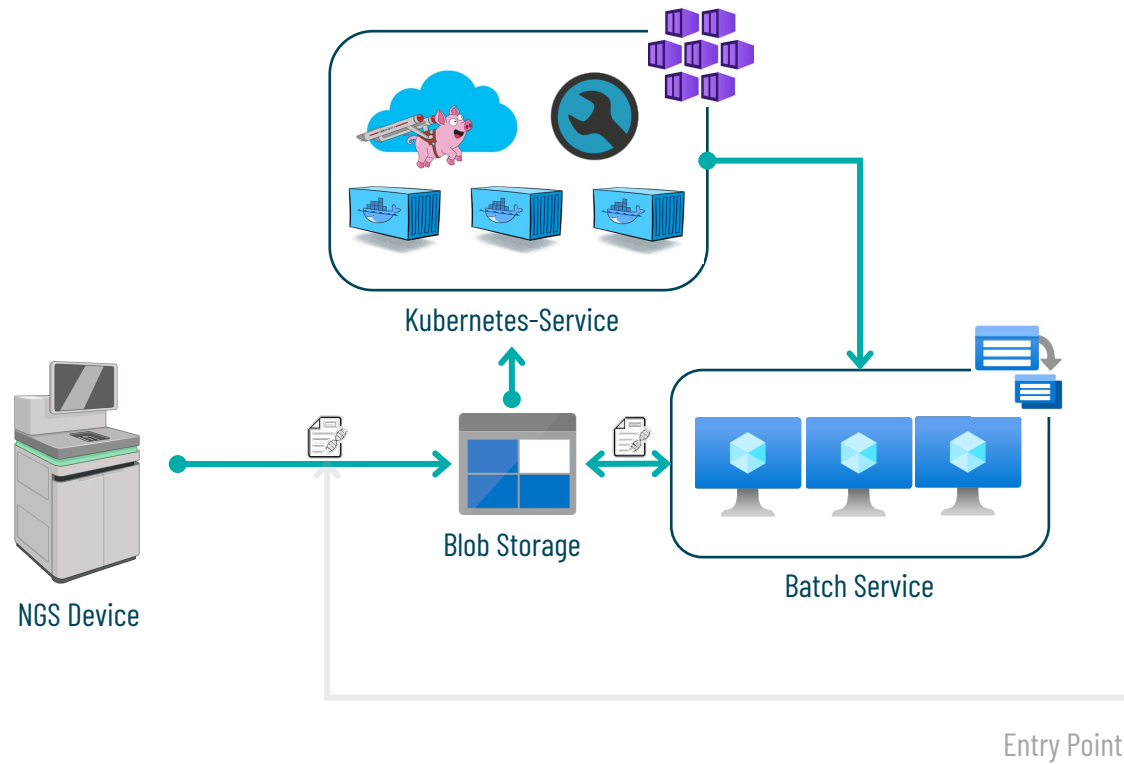


GenomAlx

Workflow & Cloud Architecture

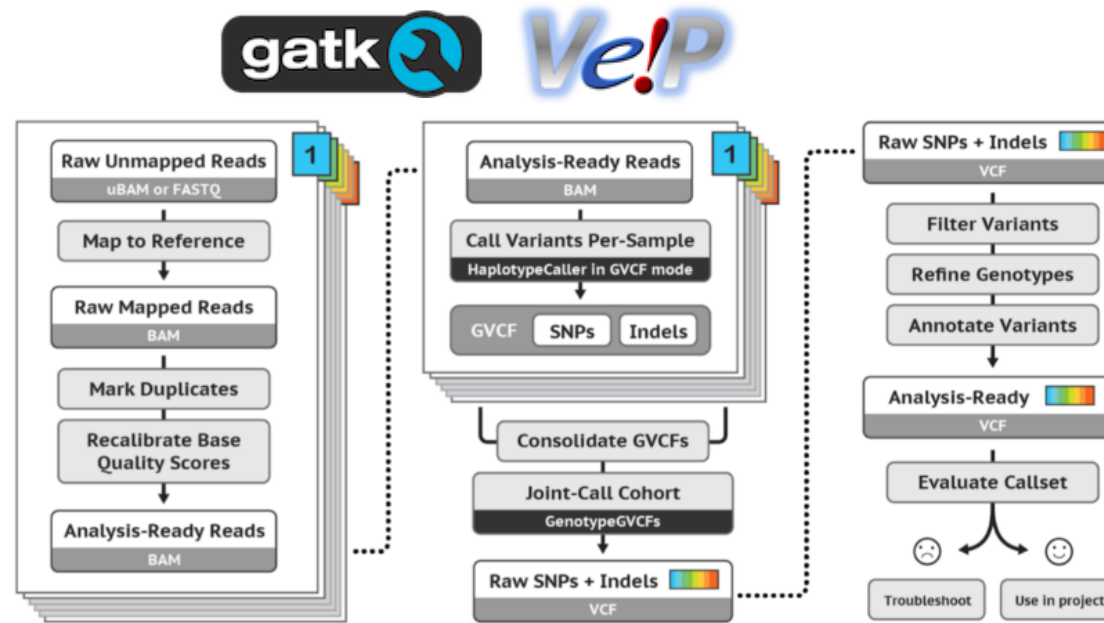
GenomAlx Workflow - Step 1

Variant Calling and Variant Effect Prediction



GenomAlx Workflow - Step 1

GATK and VEP for Variant Calling and Variant Effect Prediction

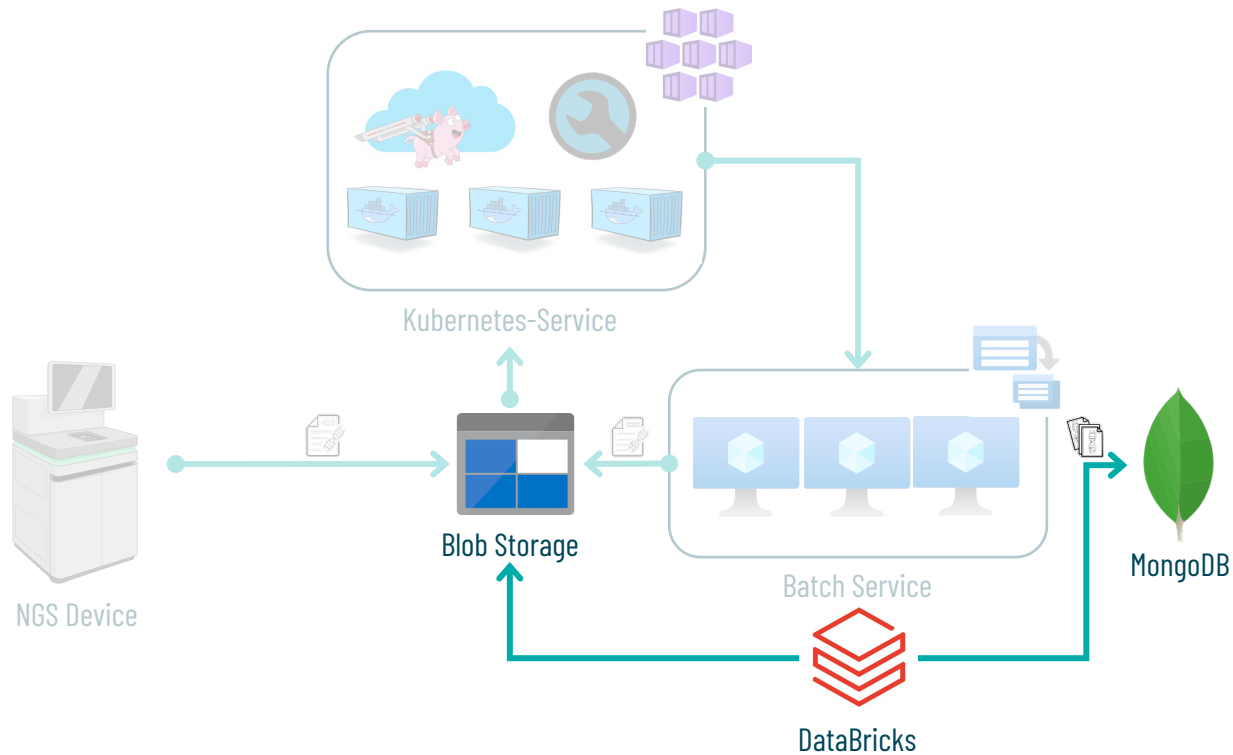


<http://www.broadinstitute.org/gatk/guide/best-practices>

```
1 Ensembl gene      1000 5000 . + . gene_id "gene1"; gene_name "GENE1";
1 Ensembl transcript 1100 4900 . + . gene_id "gene1"; transcript_id "transcript1"; gene_name "GENE1"; transcript_name "GENE1-001"; transcript_biotype "protein_coding";
1 Ensembl exon      1200 1300 . + . gene_id "gene1"; transcript_id "transcript1"; exon_number "exon1"; exon_id "GENE1-001_1";
1 Ensembl exon      1500 3000 . + . gene_id "gene1"; transcript_id "transcript1"; exon_number "exon2"; exon_id "GENE1-001_2";
1 Ensembl exon      3500 4000 . + . gene_id "gene1"; transcript_id "transcript1"; exon_number "exon3"; exon_id "GENE1-001_2";
1 Ensembl CDS       1300 3800 . + . gene_id "gene1"; transcript_id "transcript1"; exon_number "exon2"; ccds_id "CCDS0001";
```

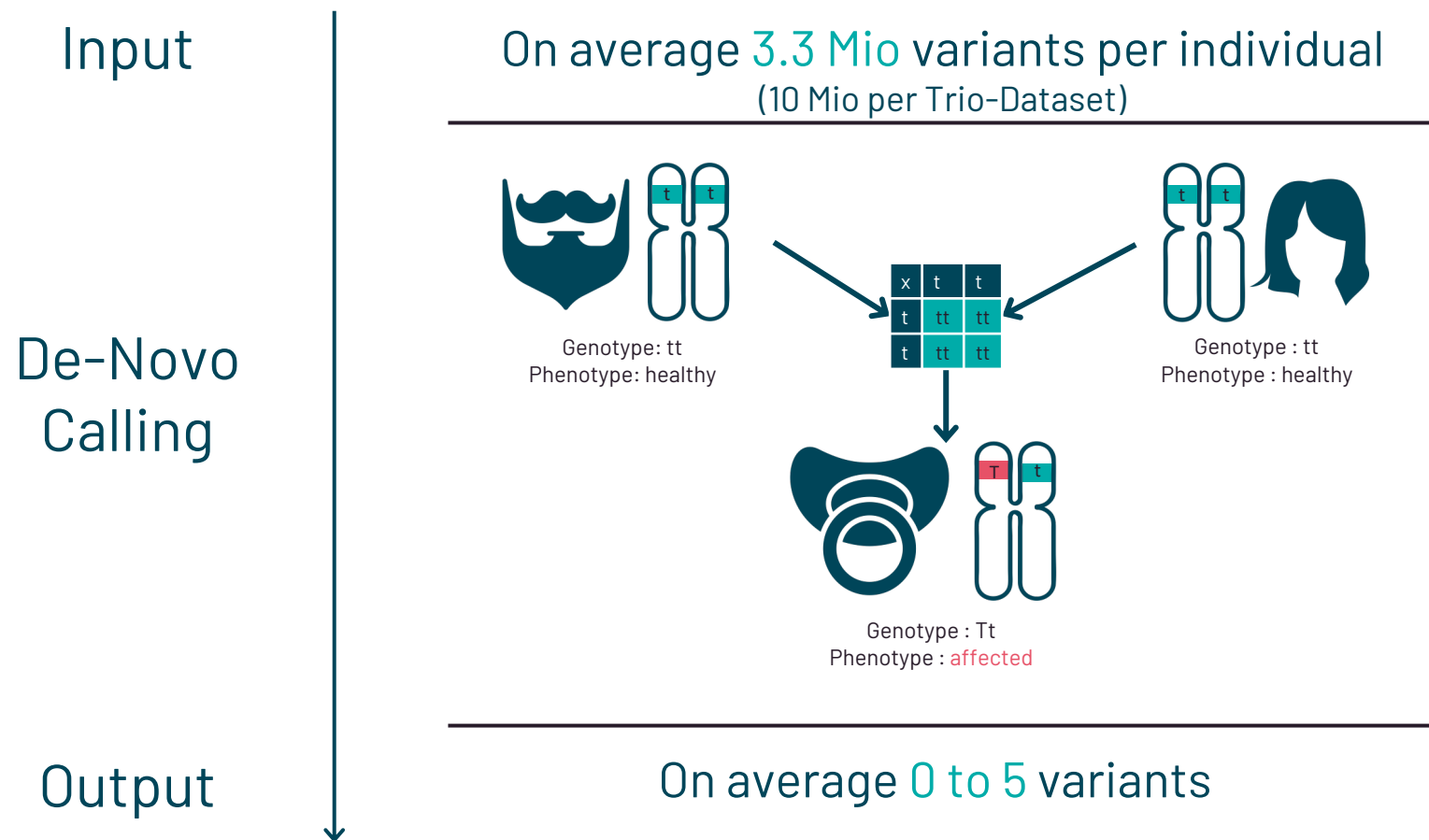
GenomAlx Workflow - Step 2

Spark based ETL-Pipelines for Constellation Calling and Variant Filtering



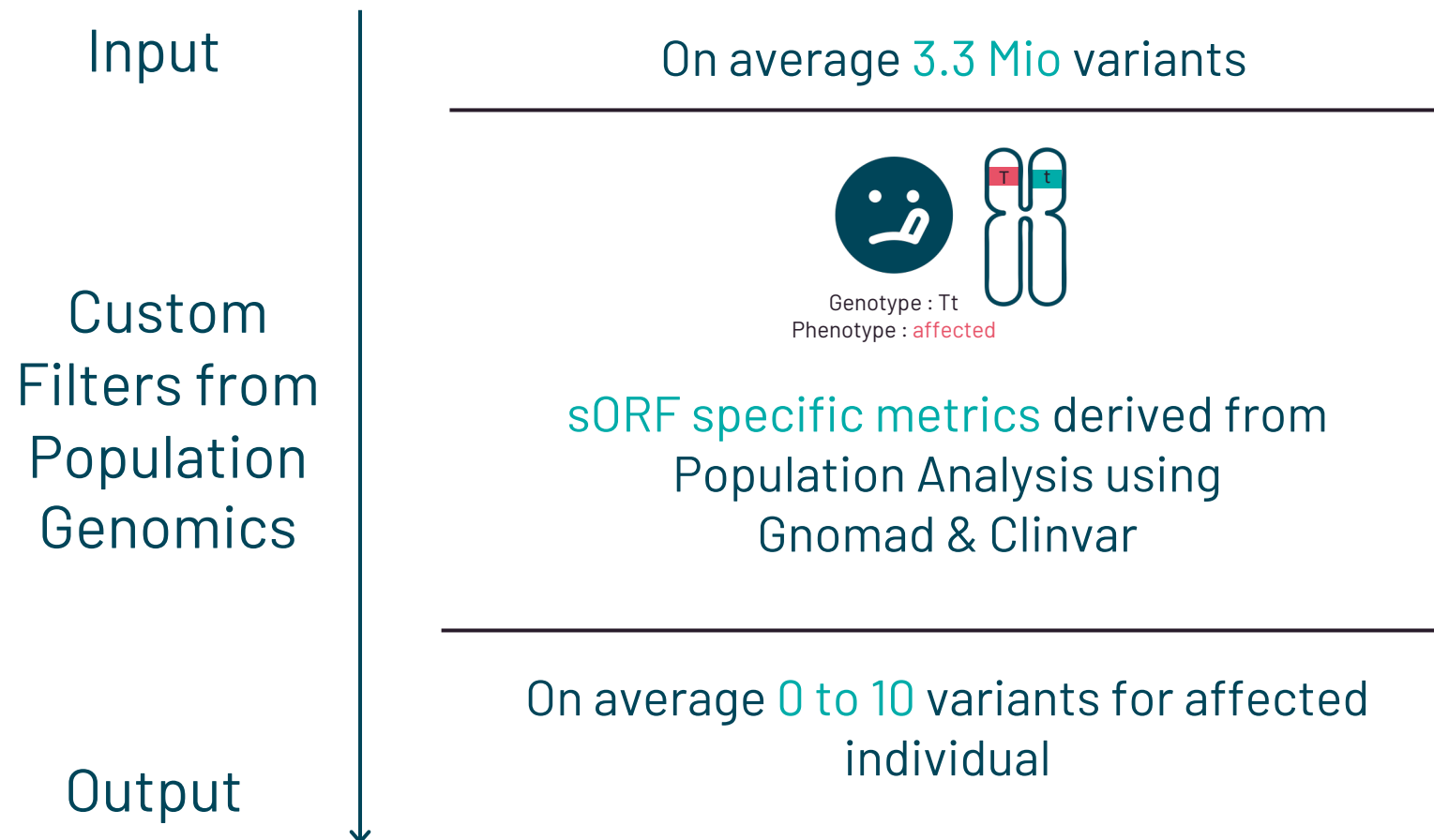
GenomAlx Workflow - Step 2

ETL-Pipelines for Constellation Calling e.g. De-Novo Calling for Trio-Datasets



GenomAlx Workflow - Step 2

ETL-Pipelines for Calling Rare Variants by means of Population Genomics



GenomAlx Workflow - Step 2

Population Genomics e.g. pLI - Score



Expected count: A depth-corrected probability prediction model that takes into account sequence context, coverage, and methylation to predict expected variant counts



Observed count: The number of unique single-nucleotide variants in each transcript/gene observed in a sample population e.g. *Gnomad* with 76.210 whole genome sequences

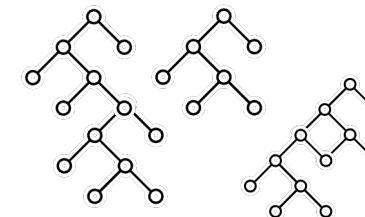


Observed/Expected ratio (O/E): When a gene has a low O/E value, it is under stronger selection for that class of variation than a gene with a higher O/E value.



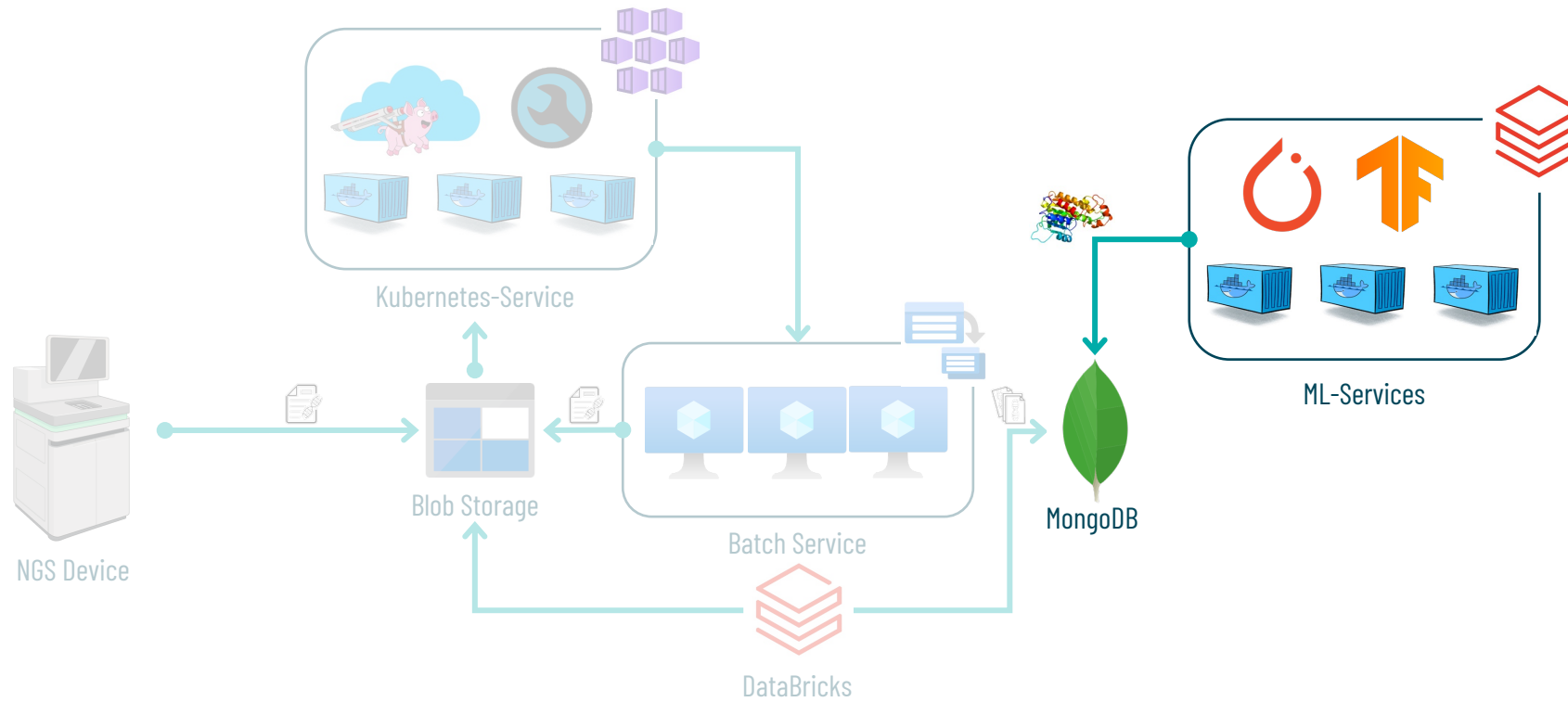
pLI-score: probability of being loss-of-function (LoF) intolerant (pLI)

XGBoost



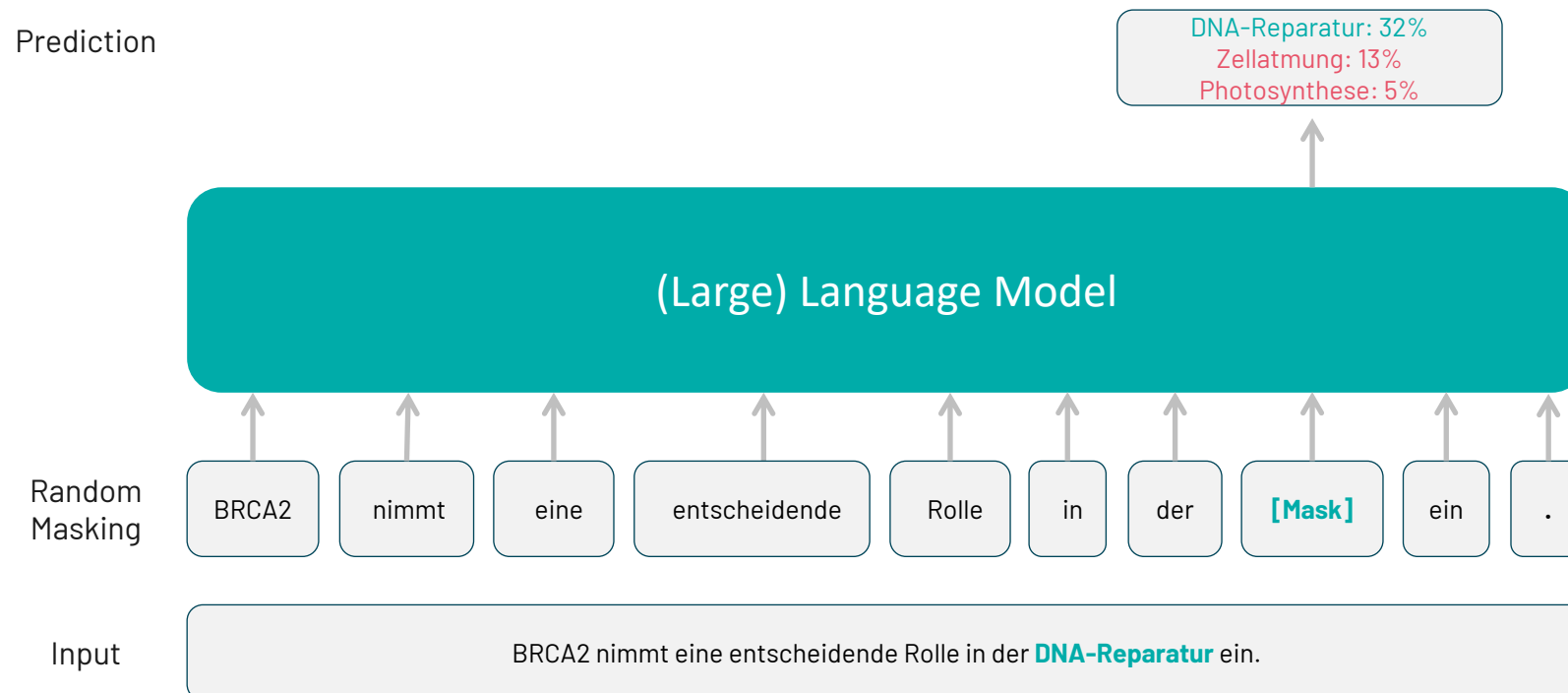
GenomAlx Workflow – Step 3

Data Enrichment using ML- Services



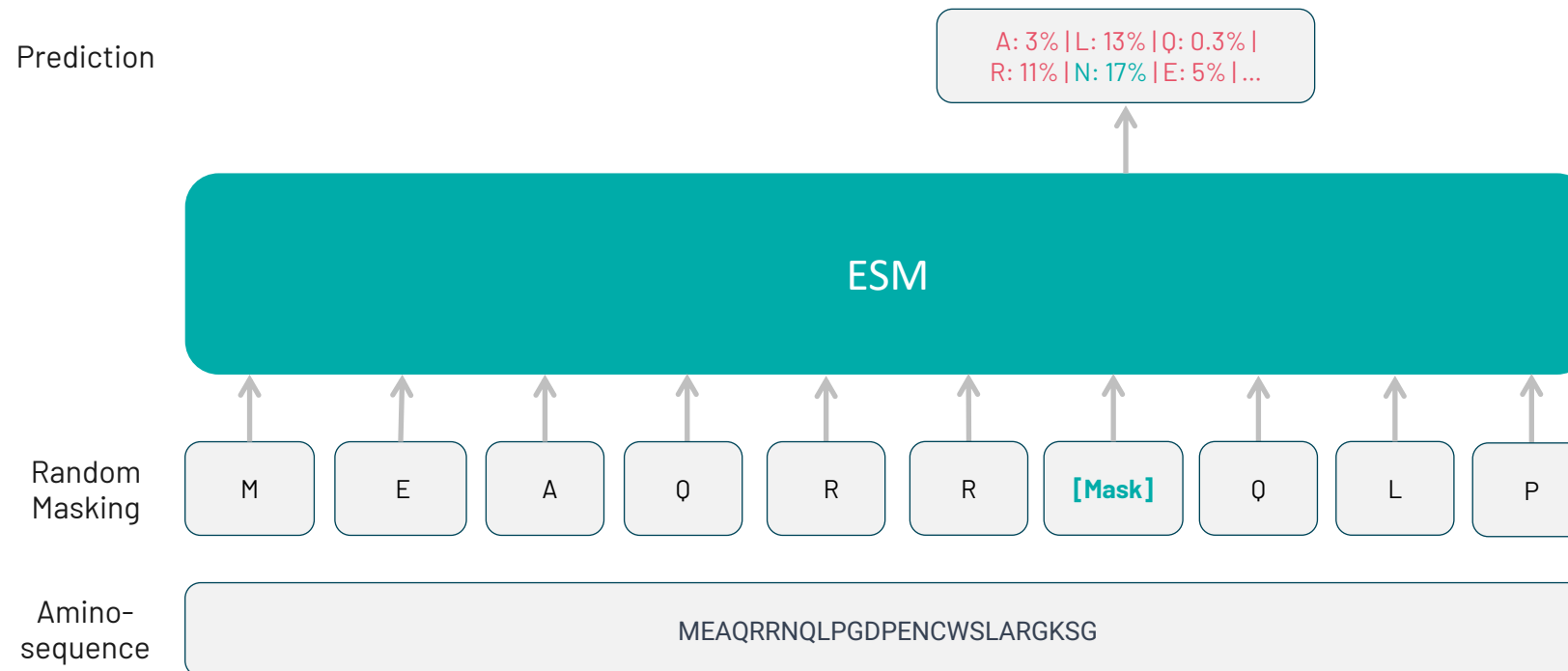
GenomAlx Workflow – Step 3

From LLMs towards Protein-Language Models via Self-Supervised Learning



GenomAlx Workflow – Step 3

From LLMs towards Protein-Language Models via Self-Supervised Learning

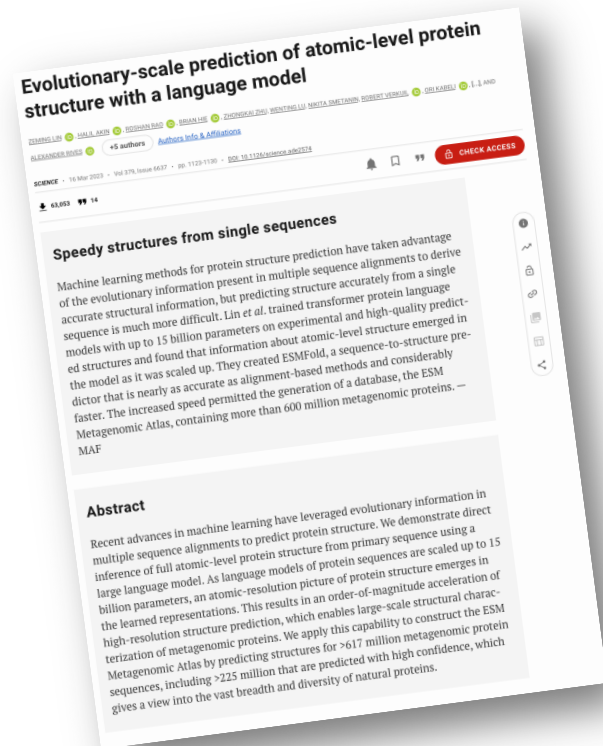


There are many similarities while comparing Aminosequences of proteins to language. Nevertheless, proteins do not have clear-cut multi-letter building blocks (such as words and sentences). They are more variable in length than sentences, and show many interactions between distant positions

[ProteinBERT: A universal deep-learning model of protein sequence and function, Brandes et al. 2022](#)

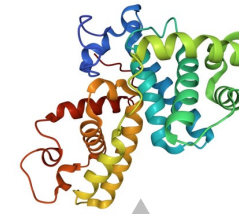
GenomAlx Workflow – Step 3

ESM-Fold: Protein Structure Prediction using Language Models

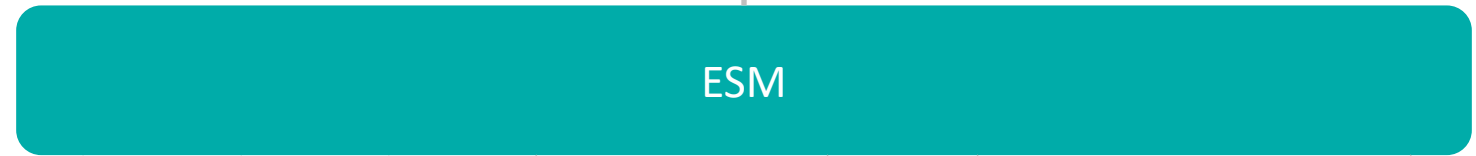
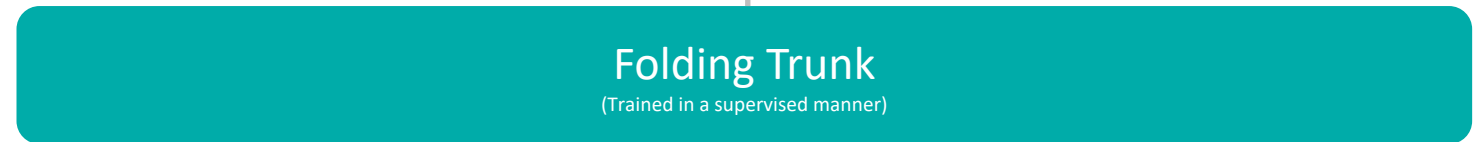


<https://www.science.org/doi/abs/10.1126/science.ade2574>

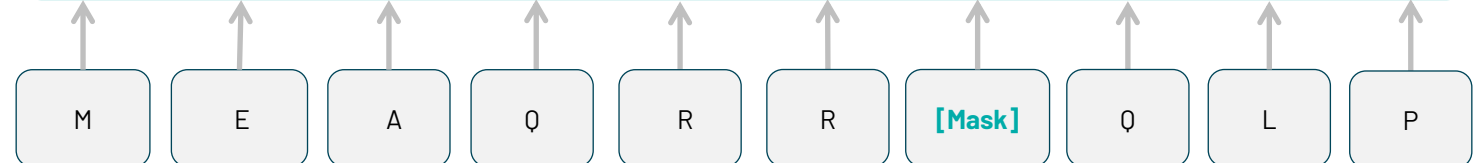
Predicted Protein Structure



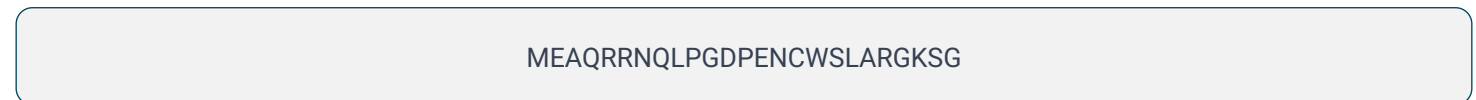
Embeddings



Random Masking



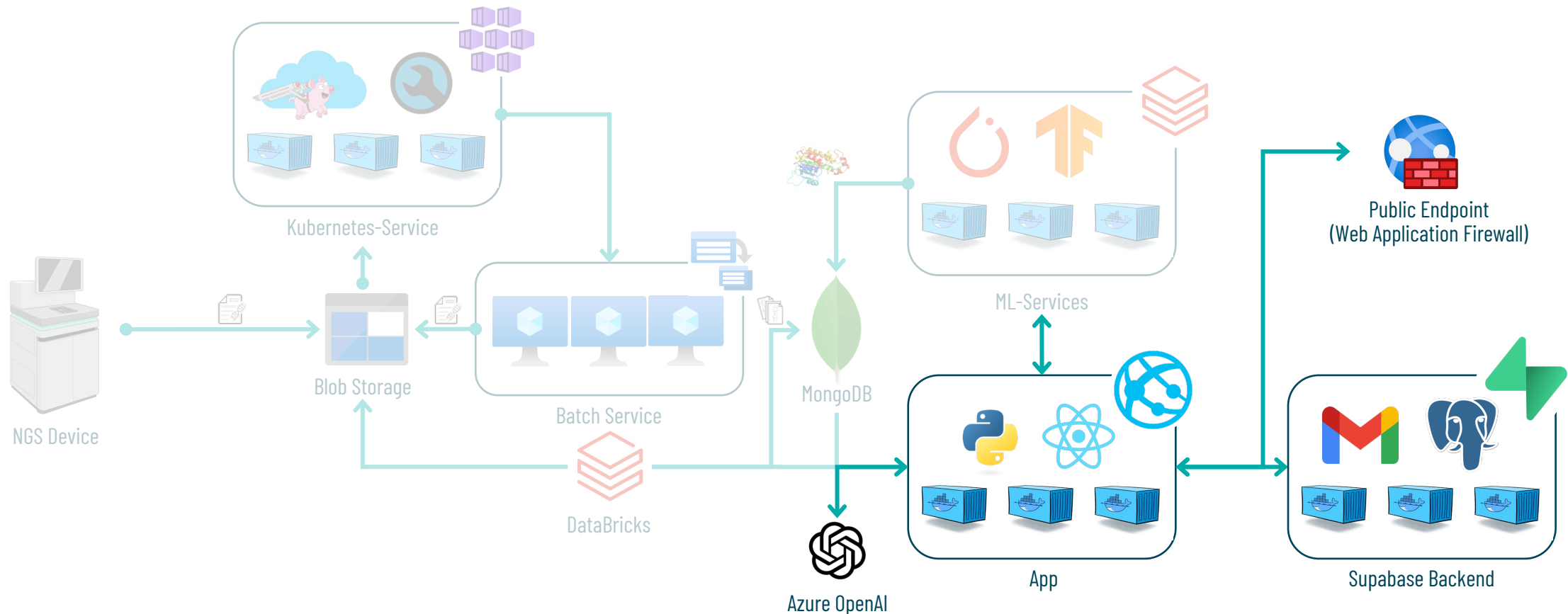
Amino-sequence



Note: Based on metrics such as perplexity ESM – Model performed worse on small aminoacid sequences (13.54) compared to larger aminoacid sequences (11.27). Since here we are focussing on proteins consisting of aminoacid sequences with a length < 150 we fully retrained ESM – Model on 12 Mio. sequences with a length < 150 AS. With that we were able to achieve better performance on small sequences with a perplexity of 12.2 for the ESM– Part of ESMFold.

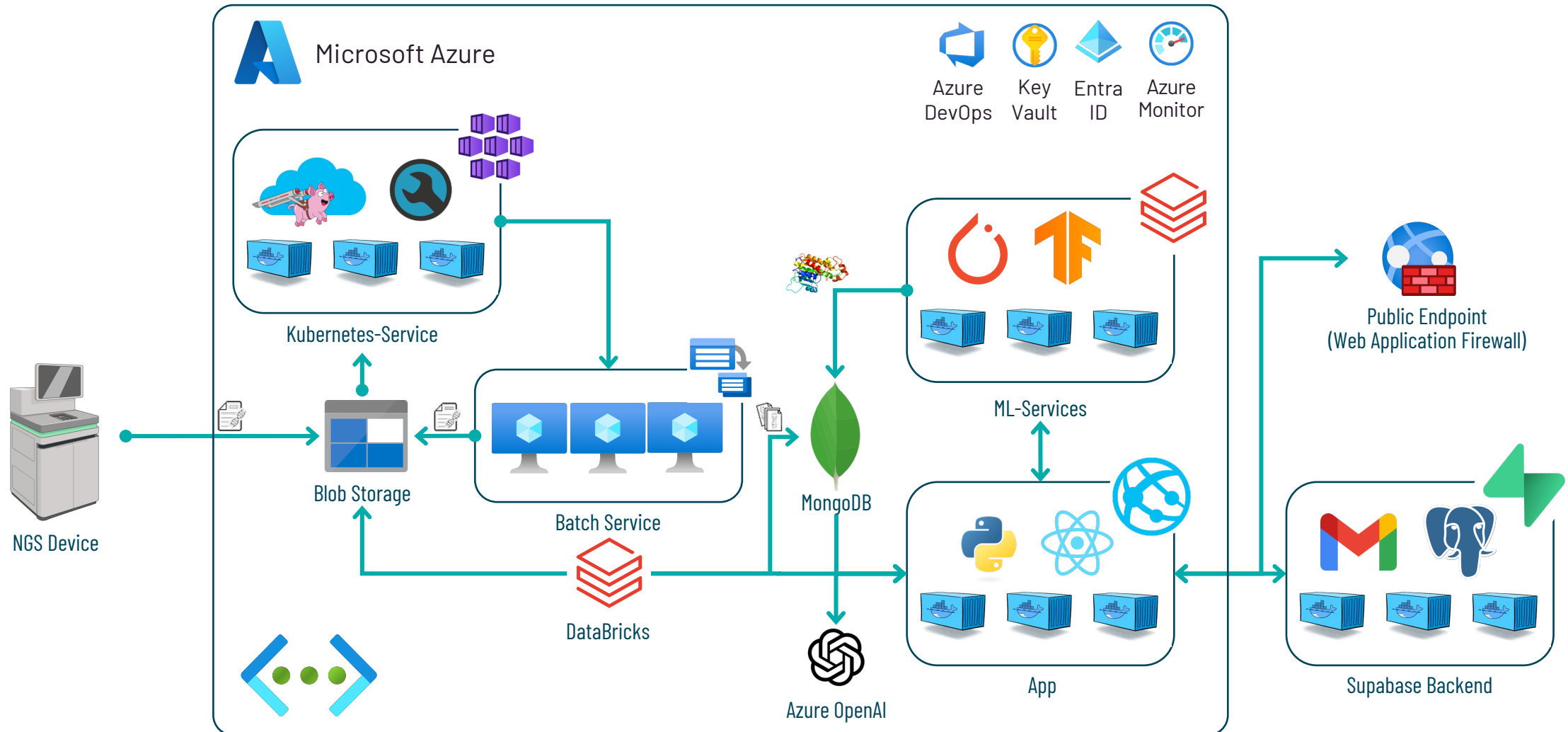
GenomAix Workflow – Step 4

Genomaix App as an Interface for Users and Human-In-The-Loop



GenomAlx Big Picture

Cloud-based Vnet secured Architecture



GenomAlx

What we achieved so far

GenomAlx

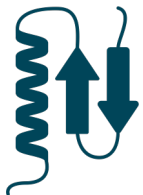
What we achieved so far



Data of 377 families processed



Data of 431 affected individuals processed



47.500 Protein Structures of sORF Variants

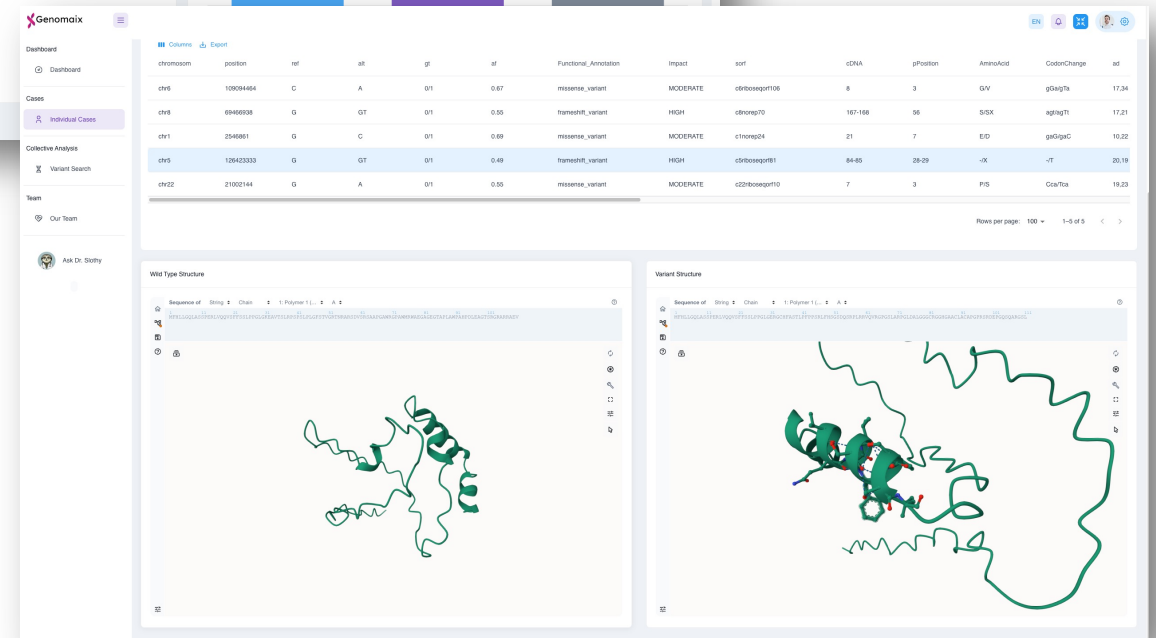
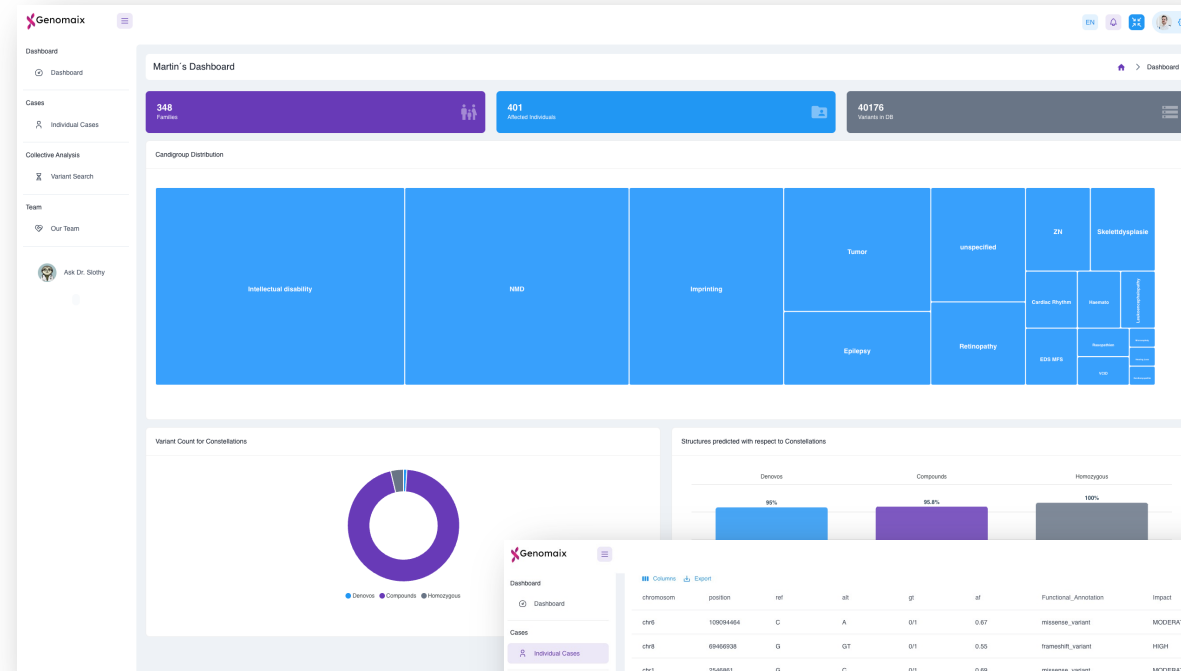


6 high interest variants identified
1 is almost certain to be disease causing

GenomAlx

von der Vision zur Applikation

(Research only)



Get in Contact!

Wir freuen uns über den Austausch.



Dr. Lars Perchalla

Direktor Data Science
scieneers GmbH

lars.perchalla@scieneers.de

Mobil +49 151 551 52 553



Prof. Dr. Ingo Kurth

Direktor Institut
Uniklinik RWTH Aachen

ikurth@ukaachen.de

Tel.: +49 241 8080178



Martin Danner

Data Scientist / PhD
scieneers GmbH / Uniklinik RWTH Aachen

martin.danner@scieneers.de

Mobil +49 151 551 52 568



Prof. Dr. Miriam Elbracht

Leitung Klinische Genomik
Uniklinik RWTH Aachen

mielbracht@ukaachen.de

Tel.: +49 241 8088013



Dr. Jeremias Krause

Assistenzarzt
Uniklinik RWTH Aachen

jkrause@ukaachen.de

Tel.: + 49 241 8087012



Dr. Matthias Begemann

Leiter NGS Diagnostikzentrum
Uniklinik RWTH Aachen

mbegemann@ukaachen.de

Tel.: + 49 241 8080036

Diskussion

