# TGAC

## The Genome Analysis Centre

*Building Excellence in Genomics and Computational Bioscience*

# The Genome Analysis Centre

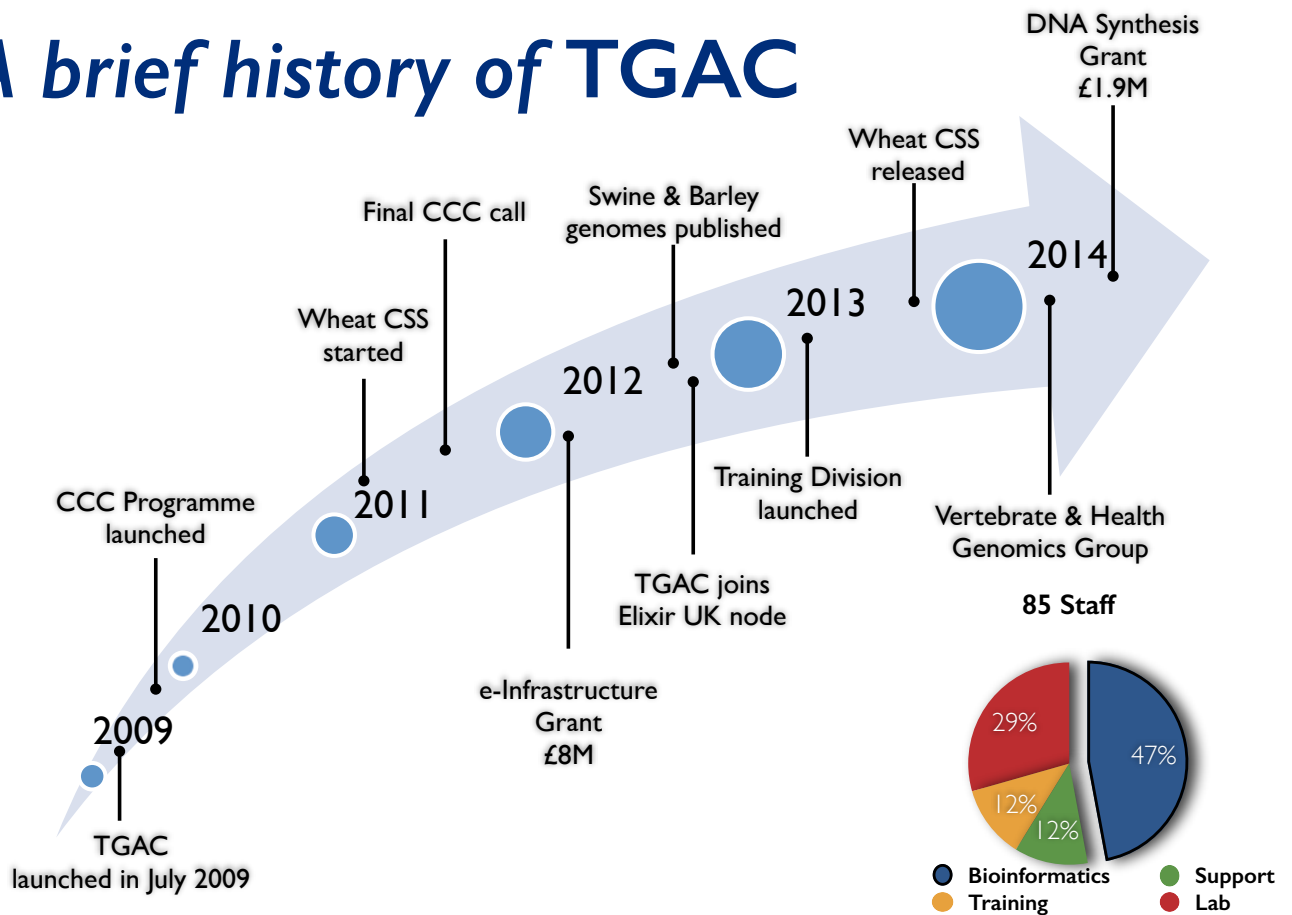### Big Data & Algorithmics in Biotechnology

## Mario Caccamo

mario.caccamo@tgac.ac.uk
@mcaccamo

**TGAC**

**The Genome Analysis Centre™**

**BBSRC**

Greater Norwich
Development
Partnership

# A brief history of TGAC

DNA Synthesis Grant £1.9M

Wheat CSS released

Swine & Barley genomes published

Final CCC call

Wheat CSS started

2014

2013

2012

CCC Programme launched

2011

2010

Training Division launched

2009

Vertebrate & Health Genomics Group

TGAC joins Elixir UK node

85 Staff

TGAC launched in July 2009

e-Infrastructure Grant £8M

47% Bioinformatics
29% Lab
12% Training
12% Support

**Bioinformatics**    **Support**
**Training**          **Lab**

TGAC ✕✕✕
The Genome Analysis Centre

# Norwich Research Park



University of East Anglia

Norfolk and Norwich University Hospital

John Innes Centre

TGAC

IFR

Sainsbury Laboratory

TGAC
The Genome Analysis Centre
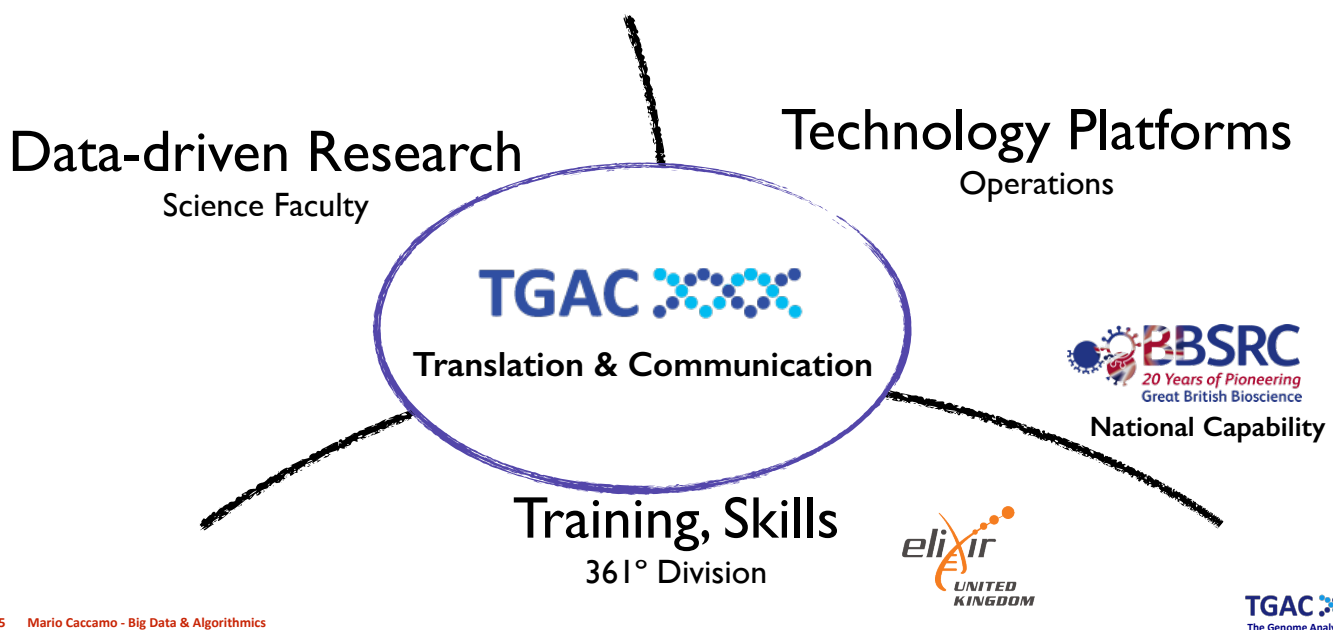
# Vision

"To be a centre of excellence in *high-end* **genomics** and **computational biology** to develop the Norwich Research Park as world leaders in bioinformatics and biotechnology."

**Data-driven Research**
Science Faculty

**Technology Platforms**
Operations

**TGAC** :✗✗✗:
**Translation & Communication**

Training, Skills
361° Division

**National Capability**

# Technology Platforms

Illumina HiSeq
**~300 Gbps/day**
**(100 times human genome)**

**Pacific Bioscience RS II
single-molecule sequencing
1of 3 in the UK**

SGI UVs
**2500 cores**
**20 Tbytes RAM**

**Lab Automation
PerkinElmer Reference Lab**

TGAC
The Genome Analysis Centre

# New Technologies



**BioNano Irys**
(Optical Mapping)



Gene Variants (Infologs)

**DNA Synthesis Grant**
**£1.9M**



**MinIon Nanopore**
**Sequencing**
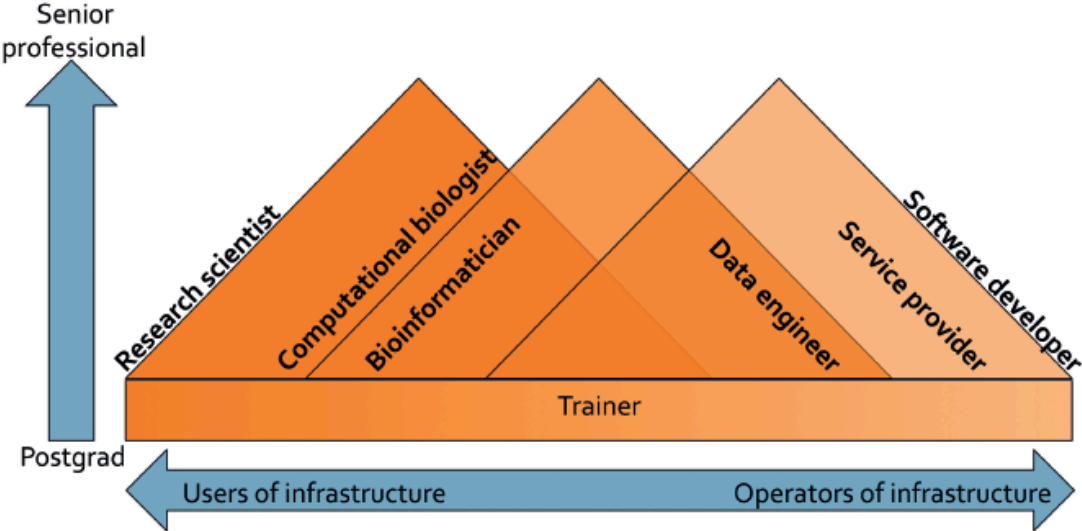(Oxford Nanopore Technologies)

# Vision

"To be a centre of excellence in *high-end* **genomics** and **computational biology** to develop the Norwich Research Park as world leaders in bioinformatics and biotechnology."

Data-driven Research
Science Faculty

Technology Platforms
Operations

**TGAC**

**Translation & Communication**

**BBSRC**
20 Years of Pioneering
Great British Bioscience

**National Capability**

Training, Skills
361° Division

elixir
UNITED KINGDOM

**TGAC**
The Genome Analysis Centre

# Elixir UK

## Bioinformatics Training!



Mario Caccamo - Big Data & Algorithmics
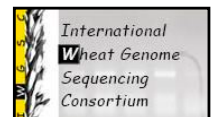
TGAC
The Genome Analysis Centre

# TransPlant



**WP12** - Implementation of resource-intensive algorithms for plant genomics data.

- To evaluate the adequacy of the current **resource-intensive** algorithms for plant genomics data.

- To build a **network of developers** with expertise in algorithms for plant genomics.

- To implement **novel algorithmic** solutions for concrete challenges in large plant genomes.
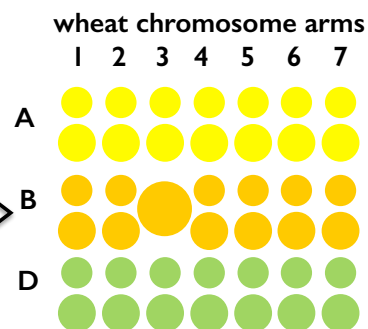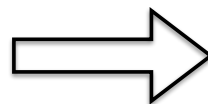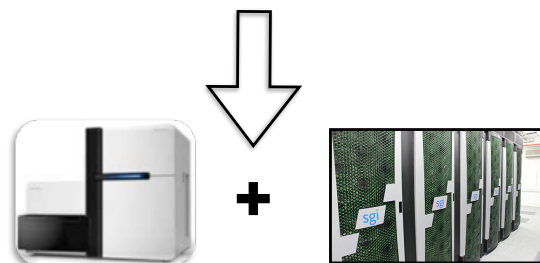
# Wheat Genome Project

**3Gb**

**Human Genome**

wheat 17Gb

wheat chromosome arms
1 2 3 4 5 6 7

A
B
D

**7.5** billion sequences
**1.5** trillion base pairs
**30** weeks of sequencing
**4.6** Terabytes
**492** hours processing time

+

10M Sequences!
First wheat
whole genome sequence
released in September 2013

*International*
*Wheat Genome*
*Sequencing*
*Consortium*

**Sarah Ayling**
**Matt Clark**

**BBSRC sLOLa Grant**
*Triticeae Genomics*
*for sustainable agriculture*
**(TGAC, JIC, EBI, RRes)**

**TGAC**
**The Genome Analysis Centre**

# From Genotype to Phenotype

# Screening for Diversity

TGAC
The Genome Analysis Centre

# Exome Captures

## Barley capture

**61.6Mb capture:**
**150,000 Morex exons**
**35,000 full length cDNAs**
**110,000 RNAseq contigs**

**IPK** GATERSLEBEN

Mascher *et al.*, The Plant Journal (2013)
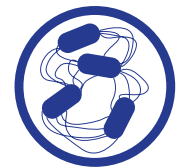
The James Hutton Institute

## Wheat capture

**84Mb capture:**
**57,000 *T. turgidum* RNAseq**
**24,000 Public wheat cDNAs**
**1,800 barley cDNAs**

**UCDAVIS**
UNIVERSITY OF CALIFORNIA

Krasileva *et al.*, Genome Biology 2013, 14:R66

John Innes Centre

TGAC
The Genome Analysis Centre

# Surveillance & Diagnostics

- **Technologies**
  - ➡ real-time data generation and analysis (Nanopore)
  - ➡ DNA extraction

- **Transmission Studies**
  - ➡ model zoonotic diseases
  - ➡ swine and horse flu transmission studies

- **Swine Flu Dynamics**
  - ➡ sLoLa (Pirbrigth, Cambridge)

- **Nornex (Ash Dieback)**
  - ➡ Genomics & genetics to tackle ash dieback
  - ➡ JIC, TSL (Norwich), Exeter, Edinburgh

TGAC
The Genome Analysis Centre

# Bioinformatics Software

| MISO | Cortex | TGAC Browser | RAMPART |
|------|--------|--------------|---------|
|  |  |  |  |

| KAT | NextClip | Bubbleparse | StatsDB |
|-----|----------|-------------|---------|
|  |  |  |  |

| RADplex | kONTAMINANT | BioJS |
|---------|-------------|-------|
|  |  |  |

**Rob Davey**

TGAC
The Genome Analysis Centre

# Bioinformatics Training

- **Launched March 2013**
  - 14 courses
  - 280 trainees

- **Elixir**
  - Elixir UK focused on Bioinformatics Training

- **Norwich Research Park**
  - Year in Industry
  - Immersive visitors training
  - TGAC Symposia

**Vicky Schneider**

# Wheat Information System

Expert Working Group within the Wheat Initiative.

Provide the wheat research community with a *single* entry point of access to genetic and genomics resources.

Promote the development of services on top of current wheat / Triticeae databases.

Authority to define guidelines for data curation, nomenclature, standards and integration.

Registry for bioinformatics tools.

TGAC
The Genome Analysis Centre

# UK Agri-Tech Strategy

- Agri-Informatics Centre for Sustainability Metrics

- Consortium led by Farming Futures (Colin Merritt) group

  ➡ NIAB, EMR, IBERS, EBI, The Cool Farm Alliance, Fera

  ➡ CAISM - Centre for Agriculture Informatics and Sustainability Metrics

**Stuart Cathpole**

TGAC
The Genome Analysis Centre