# Dark Data y los nuevos desafíos para Open Data y Open Science

Juan M. Durán

Delft University of Technology

The Netherlands

# Outline

1. Dark Data

   - Contexto institucional

   - Definición y algunas consecuencias legales

   - El Scientific Data Officer (SDO), Open Data y Open Science

   - Consequencias éticas

2. - Big Data como política de estado

**TU**Delft

# Contexto Institucional

# Dark Data and the SDO

**TU**Delft

## Dark Data as the New Challenge for Big Data Science and the Introduction of the Scientific Data Officer

Authors          Authors and affiliations

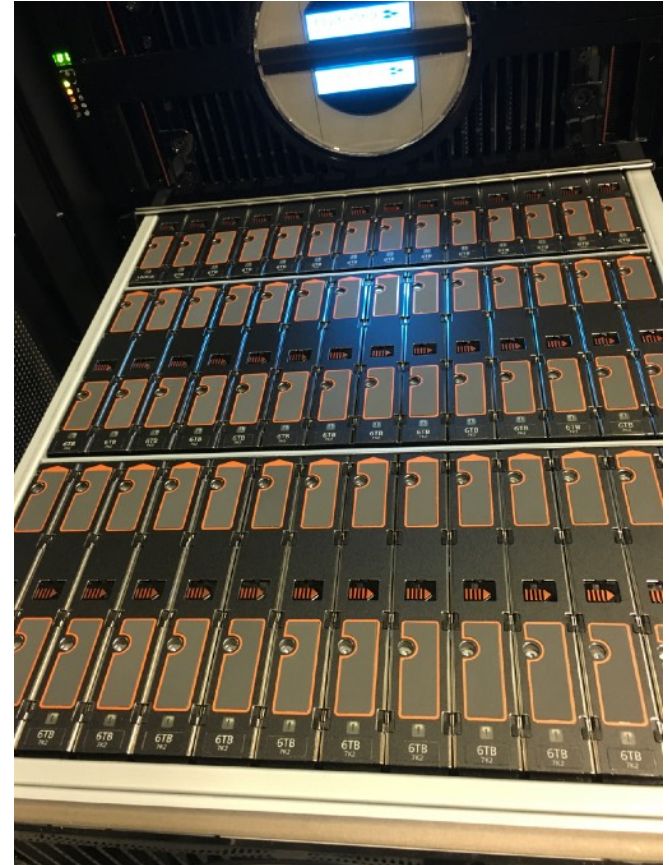Björn Schembera, Juan M. Durán ✉

## Abstract

Many studies in big data focus on the uses of data available to researchers, leaving without treatment data that is on the servers but of which researchers are unaware. We call this *dark data*, and in this article, we present and discuss it in the context of high-performance computing (HPC) facilities. To this end, we provide statistics of a major HPC facility in Europe, the High-Performance Computing Center Stuttgart (HLRS). We also propose a new position tailor-made for coping with dark data and general data management. We call it the *scientific data officer* (SDO) and we distinguish it from other standard positions in HPC facilities such as chief data officers, system administrators, and security officers. In order to understand the role of the SDO in HPC facilities, we discuss two kinds of responsibilities, namely, technical responsibilities and ethical responsibilities. While the former are intended to characterize the position, the latter raise concerns—and proposes solutions—to the control and authority that the SDO would acquire.

4

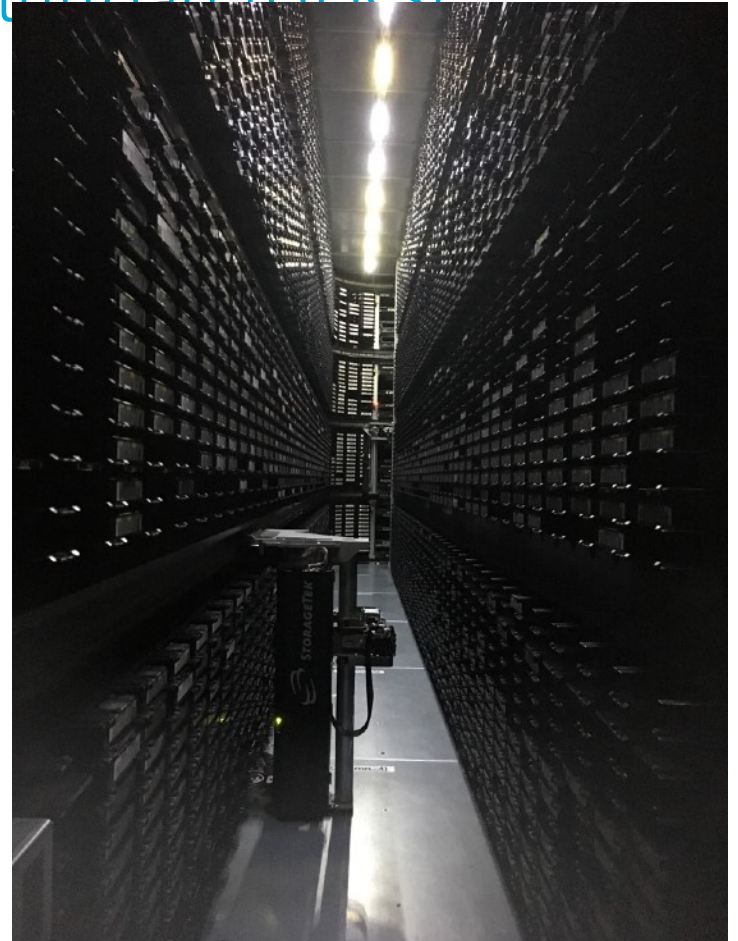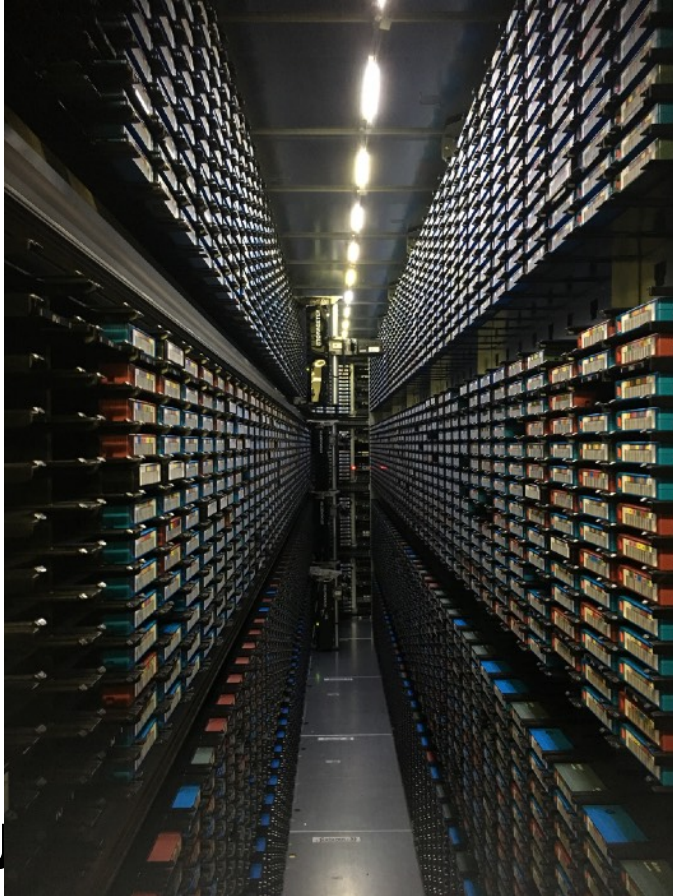# Höchstleistungsrechenzentrum Stuttgart (HLRS)



Cray XC40 Hazel-Hen

# Höchstleistungsrechenzentrum Stuttgart (HLRS)

# Höchstleistungsrechenzentrum Stuttgart (HLRS)

# Höchstleistungsrechenzentrum Stuttgart (HLRS)

- Octubre 2015

- 5,6 Peta-FLOPS per second (peak performance = 7,4 Peta-FLOPS)

- 185,088 cores

- Noviembre 2017 (3ra computadora más rápida en la UE)

- Noviembre 2015 8va más rápida del mundo. Julio 2019 cae al puesto 34.

- En 2019 tendrán una nueva "Hawk" con un costo ca. 38M€

**TU**Delft

# Höchstleistungsrechenzentrum Stuttgart (HLRS)

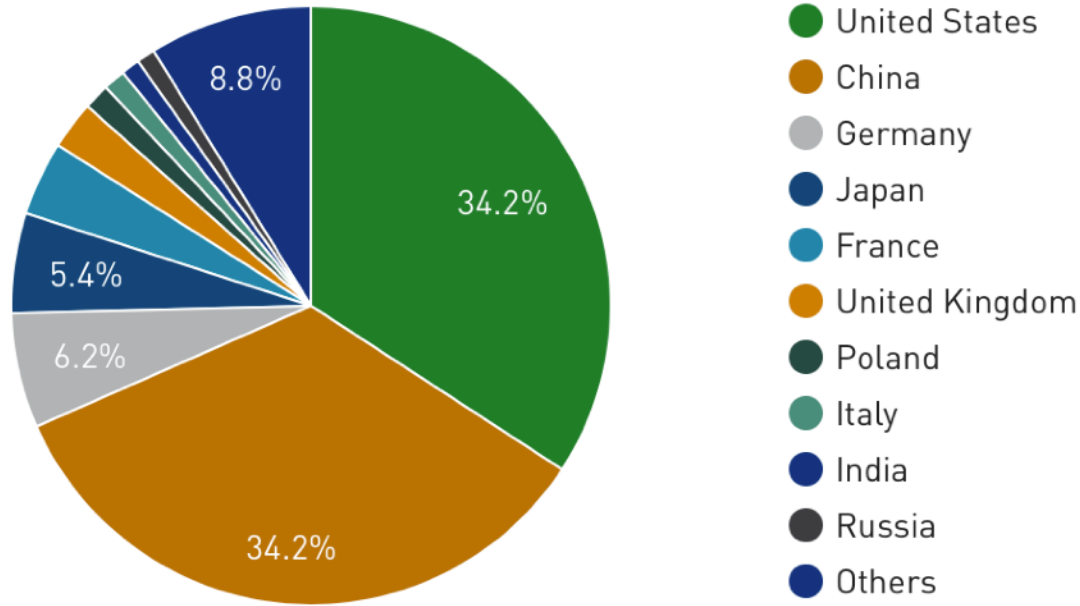| Rank | Site | System | Cores | Rmax (TFlop/s) | Rpeak (TFlop/s) | Power (kW) |
|---|---|---|---|---|---|---|
| 1 | National Super Computer Center in Guangzhou China | **Tianhe-2 (MilkyWay-2)** - TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 31S1P NUDT | 3,120,000 | 33,862.7 | 54,902.4 | 17,808 |
| | | Cray Inc. | | | | |
| 7 | Swiss National Supercomputing Centre (CSCS) Switzerland | **Piz Daint** - Cray XC30, Xeon E5-2670 8C 2.600GHz, Aries interconnect , NVIDIA K20x Cray Inc. | 115,984 | 6,271.0 | 7,788.9 | 1,754 |
| 8 | HLRS - Höchstleistungsrechenzentrum Stuttgart Germany | **Hazel Hen** - Cray XC40, Xeon E5-2680v3 12C 2.5GHz, Aries interconnect Cray Inc. | 185,088 | 5,640.2 | 7,403.5 | 3,615 |
| 9 | King Abdullah University of Science and Technology Saudi Arabia | **Shaheen II** - Cray XC40, Xeon E5-2698v3 16C 2.3GHz, Aries interconnect Cray Inc. | 196,608 | 5,537.0 | 7,235.2 | 2,834 |

https://www.top500.org Nov 2015

# Höchstleistungsrechenzentrum Stuttgart (HLRS)

| Rank | Site | System | Cores | Rmax (TFlop/s) | Rpeak (TFlop/s) | Power (kW) |
|---|---|---|---|---|---|---|
| 1 | National Supercomputing Center in Wuxi China | **Sunway TaihuLight** - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway NRCPC | 10,649,600 | 93,014.6 | 125,435.9 | 15,371 |
| 2 | National Super Computer Center in Guangzhou China | **Tianhe-2 (MilkyWay-2)** - TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 31S1P NUDT | 3,120,000 | 33,862.7 | 54,902.4 | 17,808 |
| | | HPE/SGI | | | | |
| 14 | HLRS - Höchstleistungsrechenzentrum Stuttgart Germany | **Hazel Hen** - Cray XC40, Xeon E5-2680v3 12C 2.5GHz, Aries interconnect Cray Inc. | 185,088 | 5,640.2 | 7,403.5 | 3,615 |
| 15 | King Abdullah University of Science and Technology Saudi Arabia | **Shaheen II** - Cray XC40, Xeon E5-2698v3 16C 2.3GHz, Aries interconnect Cray Inc. | 196,608 | 5,537.0 | 7,235.2 | 2,834 |

# Distribución de HPC

**Countries System Share**



- United States — 34.2%
- China — 34.2%
- Germany — 6.2%
- Japan — 5.4%
- France
- United Kingdom
- Poland
- Italy
- India
- Russia
- Others — 8.8%

https://www.top500.org Nov 2016

# AméricaLatina (bueno, Brasil…)

| | | | | | | |
|---|---|---|---|---|---|---|
| 142 | Petróleo Brasileiro S.A<br>Brazil | **Fênix** - SYS-1029GQ-TRT, Xeon Gold 5122 4C<br>3.6GHz, Infiniband EDR, NVIDIA Tesla V100<br>Bull, Atos Group | 48,384 | 1,836.0 | 4,297.4 | 287 |
| 418 | Cloud Provider<br>Brazil | **BC1** - Lenovo C1040, Xeon E5-2673v4 20C 2.3GHz,<br>40G Ethernet<br>Lenovo | 38,400 | 1,123.2 | 1,413.1 | |
| 419 | Cloud Provider | Lenovo C1040, Xeon E5-2673v4 20C 2.3GHz, 40G | 38,400 | 1,123.2 | 1,413.1 | |
| 429 | Software Company (M)<br>Brazil | **BC2** - Lenovo C1040, Xeon E5-2673v4 20C 2.3GHz,<br>40G Ethernet<br>Lenovo | 38,400 | 1,123.2 | 1,413.1 | |

# Dark Data

# Dark Data - una definición posible

"Dark data" is [data] not carefully indexed and stored so it becomes nearly invisible to scientists and other potential users and therefore is more likely to remain under-utilised and eventually lost [..] the type of data that exists only in the bottom left-hand desk drawer of scientists on some media that is quickly ageing and soon will be unreadable by commonly available devices" (Heidorn 2008, 281)

**TU**Delft

# Dark Data - una definición posible

"Dark data" is [data] not carefully **indexed** and stored so it becomes nearly **invisible** to scientists and other potential users and therefore is more likely to remain **under-utilised** and eventually **lost** [..] the type of data that exists only in the bottom left-hand desk drawer of scientists on some media that is quickly ageing and soon will be **unreadable** by commonly available devices" (Heidorn 2008, 281)

**TU**Delft

# Dark Data - Metadata

Para nosotros"dark data" también son:

- Datos que no tienen meta-datos (i.e., no han sido etiquetados), que están mal organizados (tal vez por limitaciones en el filesystem), que están dispersos.
  - Son difícil de contabilizar: definiciones y mediciones son necesarias

| Dec-2017 | 668.03 TB / 19,663 TB (3,40 % DD) |
|----------|-----------------------------------|
| Mar-2019 | 1,021.40 TB / 21,240 TB (4,81 % DD) |

# Dark Data - Usuarios

Para nosotros "dark data" también un origen en los usuarios del sistema:

- Datos guardados que persisten en el sistema de archivos por usuarios inactivos o des-registrados
  - Cómo determinar que un usuario no está más activo?
    - Tiempo inactivo?
    - Cantidad de datos transferidos?

# Por qué hay presencia de Dark Data?

- Mantener los datos y los metadatos no motiva, ni profesionalmente, ni económicamente, ni socialmente. En otras palabras, por qué utilizar "tiempo precioso" para organizar mis datos (más allá de lo necesario para que yo pueda acceder a esos datos)

- Centros de HPC adoptan nueva infraestructura, nuevos in-house sistemas de archivos, contribuye a los dark data.

**TU**Delft

# Implicaciones del Dark Data

- Recursos:
    - Dark Data desperdicia recursos (almacenamiento, tiempo de cómputo, mantenimiento)
    - 500TB de Dark Data cuestan ca. $ 10,000 dependiendo de la tecnología de almacenamiento de datos.
    - Se invierte dinero en educar a un ingeniero/científico.

- En varios sentidos el Dark Data viola regulaciones del GDPR
    - e.g., cuando los datos contienen información personal
    - e.g., cuando los datos tienen derechos de propiedad

- Dark Data no es FAIR: **F**indable, **A**ccessible, **I**nter-operable, **R**eusable (por definición!)

**TU**Delft

# Scientific Data Officer (SDO)

# Scientific Data Officer: necesidad & responsabilidades

- Etiquetar (o dar soporte en el etiquetado) de datos con algún formato estandarizado de metadatos de acuerdo con FAIR (inexistente!)
- Estandarizar los datos y metadatos para todos los proyectos en un centro de HPC
- Reduce Dark Data
- Chequeos periódicos del inventario de datos (y de Dark Data)
- Administración del sistema de datos
- Autoridad sobre el problema de decidir sobre usuarios inactivos/des-registrados.

**TU**Delft

# Scientific Data Officer: responsabilidades

- Stewardship/curadora de datos

- Multiplicadora de conocimiento y su transferencia (no se necesita más tener "el contacto" o conocer "al juez")

- Rol mediadora (e.g., comparte información y datos intra e inter-institucional)
  - Una red de SDOs?

- Dar apertura a los datos (cumplir con los ideales de OpenScience and OpenData) y cumplir con FAIR

- La SDO debe observar y hacer respetar reglas de buena práctica científica (incluido la autoridad de sancionar y accionar legalmente?).

**TU**Delft

# Scientific Data Officer: Problemas éticos

- Se presentan varias consideraciones sobre el comportamiento de las SDOs:
  - Malversación de los resultados de una investigación (embezzlement/ missappropriation)
    - Manipulation y selección de datos
  - Interferencia en la producción de resultados
  - Obstrucción en la investigación científica
    - Priorizar el acceso al sistema de archivos a ciertos grupos en desventaja de otros
  - No cumplimientos con códigos de conducta y éticos
  - Gatekeeping: determina quién continua con una pertenencia a la institución

**TU**Delft

# Big Data como política de estado

# Making Dark Data FAIR



EUROPEAN OPEN SCIENCE CLOUD - EOSC
Funding Opportunities with
Co-Creation Budget

EOSCsecretariat.eu has received funding from the European Union's Horizon Programme call H2020-INFRAEOSC-2018-4, grant Agreement number 831644

# SoBigDataPlusPlus



SoBigData is the European Research Infrastructure for **Big Data** and Social Mining. From data to knowledge, investigating stories **ethically**, paying attention to citizens privacy.

http://www.sobigdata.eu/index

# SoBigDataPlusPlus



https://sobigdata.d4science.org/group/sobigdata-gateway

# HumaineAI

# HumaineAI

We are designing the principles for a new science that will make artificial intelligence based on European values and closer to Europeans.

This new approach works toward AI systems that augment and empower all Humans by understanding us, our society and the world around us.

**TU**Delft

# HumaineAI

- Valores Europeos?

- China y EEUU:

  - China: Estado policíaco de vigilancia (surveillance state). Para 2018 China ya tenía una red de vigilancia de más de 170 millions de cámaras CCTV, con aprox. 400 million de nuevas cámaras a instalarse en los próximos tres años (el uso es el "facial recognition technology") - Financiamieno es "virtualmene infinito"

**TU**Delft

# HumaineAI

- Valores Europeos?

- China y EEUU:

  - EEUU: El gran motor en Big Data es la industria, con sus valores capitalistas y con fines comerciales y obtención de ganancias. Esto, en principio, podría afectar libertades y derechos adquiridos, como el derecho a la privacidad, a la protección del estado, al acceso a tecnologías, etc.

**TU**Delft

# HumaineAI

- Valores Europeos
  - Estado presente regulando la industria pero a su vez monitoreado por ONG y agencias independientes que garanticen los derechos de todos los ciudadanos.
  - GDPR

**TU**Delft

# HumaineAI



http://videolectures.net/humaneai_duran_interview/

33

Más información →

# Muchas Gracias

Juan M. Durán - TU Delft

j.m.duran@tudelft.nl

juanmduran.net

**TU**Delft

## Dark Data as the New Challenge for Big Data Science and the Introduction of the Scientific Data Officer

Björn Schembera[1] · Juan M. Durán[2]

**Abstract**
Many studies in big data focus on the uses of data available to researchers, leaving without treatment data that is on the servers but of which researchers are unaware. We call this *dark data*, and in this article, we present and discuss it in the context of high-performance computing (HPC) facilities. To this end, we provide statistics of a major HPC facility in Europe, the High-Performance Computing Center Stuttgart (HLRS). We also propose a new position tailor-made for coping with dark data and general data management. We call it the *scientific data officer* (SDO) and we distinguish it from other standard positions in HPC facilities such as chief data officers, system administrators, and security officers. In order to understand the role of the SDO in HPC facilities, we discuss two kinds of responsibilities, namely, technical responsibilities and ethical responsibilities. While the former are intended to characterize the position, the latter raise concerns—and proposes solutions—to the control and authority that the SDO would acquire.

**Keywords** Research data management · High-performance computing · Dark data · Big data · Computer simulations · Scientific data officer · Data curation

Both authors contributed equally to all sections of the paper.

✉ Juan M. Durán
j.m.duran@tudelft.nl

Björn Schembera
schembera@hlrs.de

[1] High-Performance Computing Center Stuttgart, University of Stuttgart, Nobelstr. 19, 70569 Stuttgart, Germany

[2] Faculty of Technology, Policy and Management, Delft University of Technology, Jaffalaan 5, 2628 BX Delft, Netherlands

Springer