# Pediatric Cancer Research Data Commons

## Samuel Volchenboum, MD, PhD
## September, 2017

THE UNIVERSITY OF CHICAGO MEDICINE & BIOLOGICAL SCIENCES | CENTER FOR RESEARCH INFORMATICS
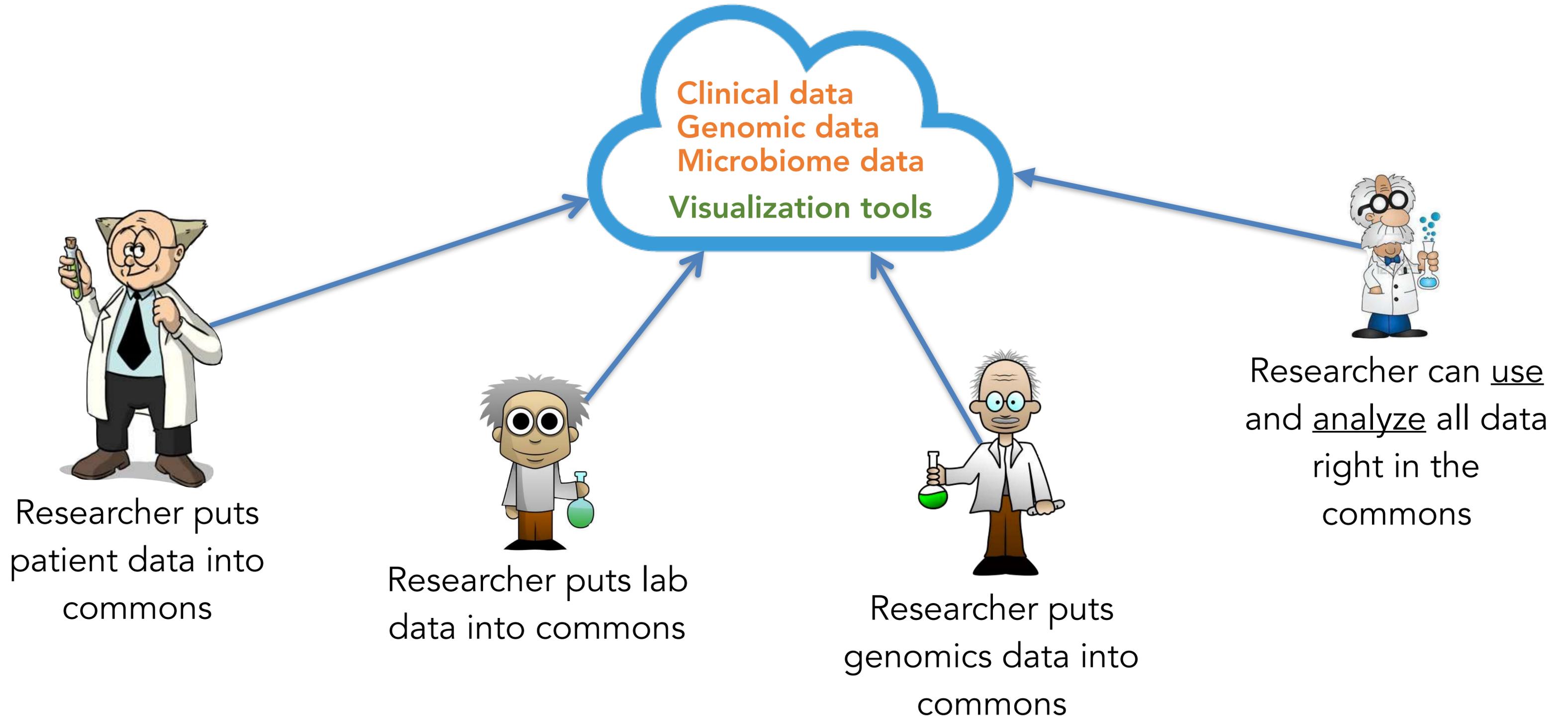
Most research is performed in silos

# A commons is a place to store and use data



Clinical data
Genomic data
Microbiome data
Visualization tools

Researcher puts patient data into commons

Researcher puts lab data into commons

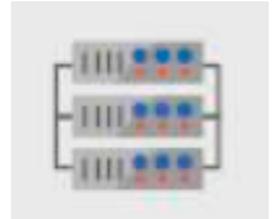Researcher puts genomics data into commons

Researcher can use and analyze all data right in the commons

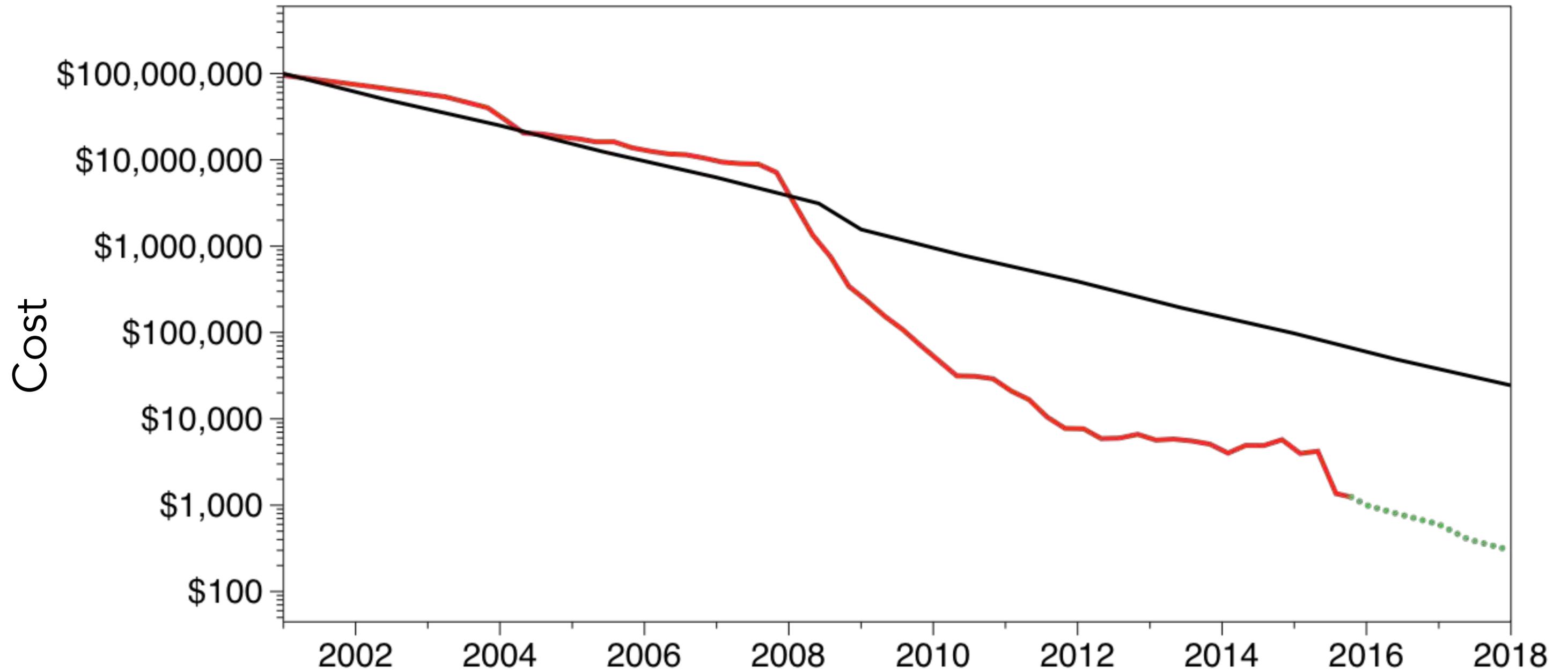# Why do we need data commons?

Too much data to store

Takes too long to transfer

Too expensive to analyze

Lack of data standardization

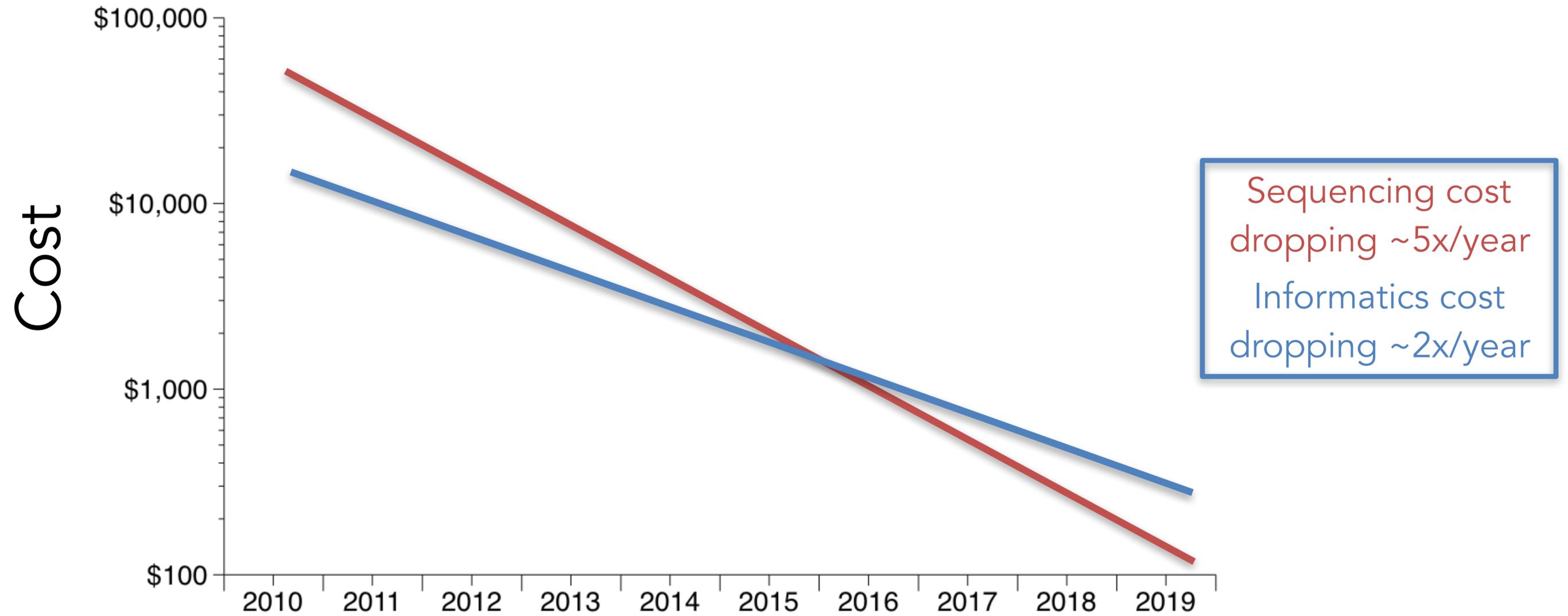# Falling costs of DNA sequencing

# Informatics is now the bottleneck



Cost

$100,000
$10,000
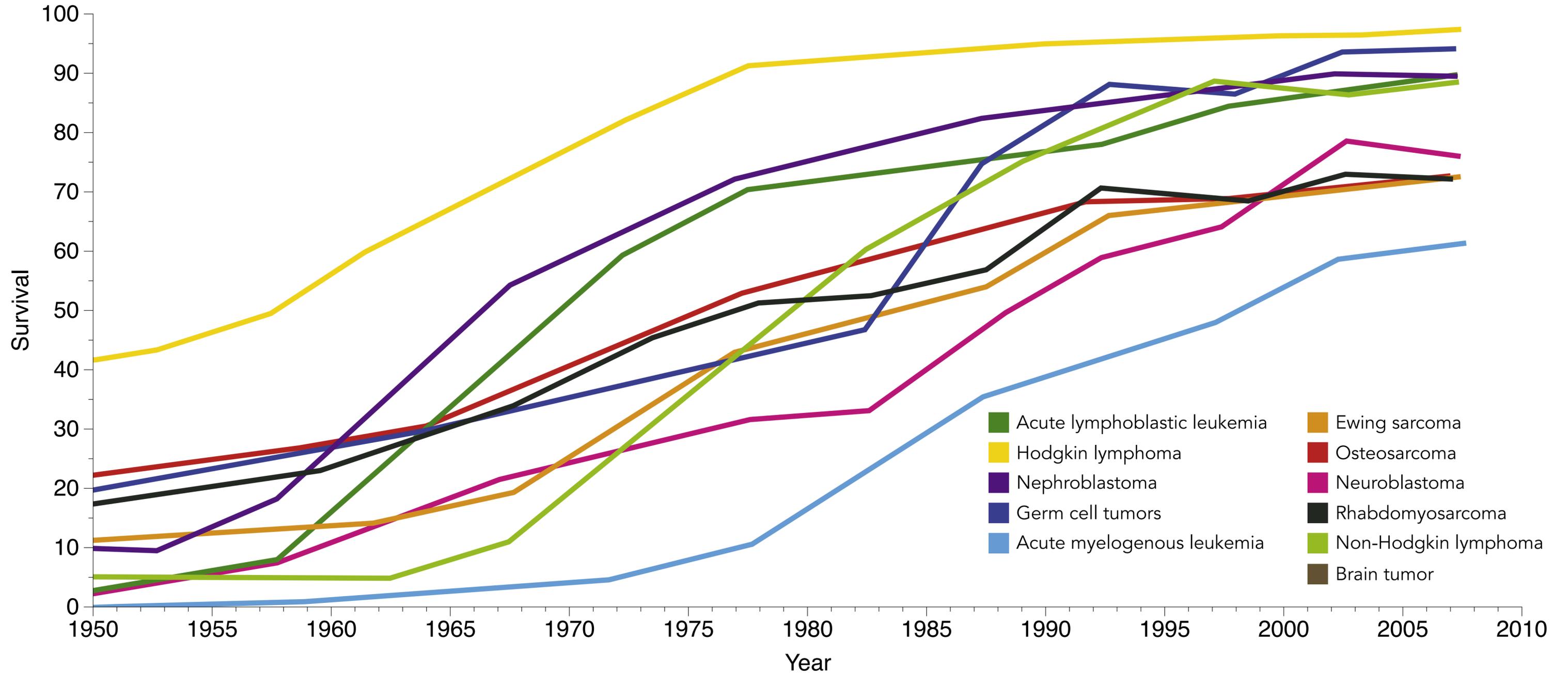$1,000
$100

2010 2011 2012 2013 2014 2015 2016 2017 2018 2019

Sequencing cost dropping ~5x/year
Informatics cost dropping ~2x/year

businessinsider.com

THE UNIVERSITY OF CHICAGO MEDICINE & BIOLOGICAL SCIENCES

CENTER FOR RESEARCH INFORMATICS

# The importance of standards

Pediatric rhabdomyosarcoma

| | | Current | |
|---|---|---|---|
| | Code | Description | |
| **SEX** | 1 | Male | |
| | 2 | Female | |
| **HISTOLOGY** | 1 | Alveolar rhabdomyosarcoma | |
| | 2 | Embryonal rhabdomyosarcoma | |
| | 3 | Botryoid rhabdomyosarcoma | |
| | 4 | Not otherwise specified | |
| | 5 | Undifferentiated sarcoma | |
| | 6 | Sarcoma, not classifiable | |
| | 7 | Spindle cell sarcoma | |
| | 8 | Ectomesenchymoma | |
| | 9 | Other | |
| | 10 | Mixed rhabdomyosarcoma | |
| | 99 | Unknown | |

F - Findable
A - Accessible
I - Interoperable
R - Reusable

# Survival - Pediatric cancer

# Why do we need <u>pediatric</u> cancer data commons?

## Adult cancers

| All | 1,688,780 |
|---|---|
| Oral | 49,670 |
| GI | 310,440 |
| Lung | 222,500 |
| Skin | 95,360 |
| Breast | 255,180 |
| Ovary | 22,440 |
| Prostate | 161,360 |
| Urinary | 146,650 |
| Lymphoma | 80,500 |
| Myeloma | 30,280 |
| Leukemia | 62,130 |

Source: cancer.org - 2017

## Pediatric cancers

| All | 15,190 |
|---|---|
| Bone | 787 |
| Brain | 2,653 |
| Hodgkin lymphoma | 973 |
| Kidney | 608 |
| Acute lymphoblastic leukemia | 2,627 |
| Acute myelogenous leukemia | 642 |
| Non-Hodgkin lymphoma | 1,116 |
| Soft-tissue | 869 |
| Other | 4,585 |

Source: CDC - 2014

Studying a rare disease like pediatric cancer is difficult.
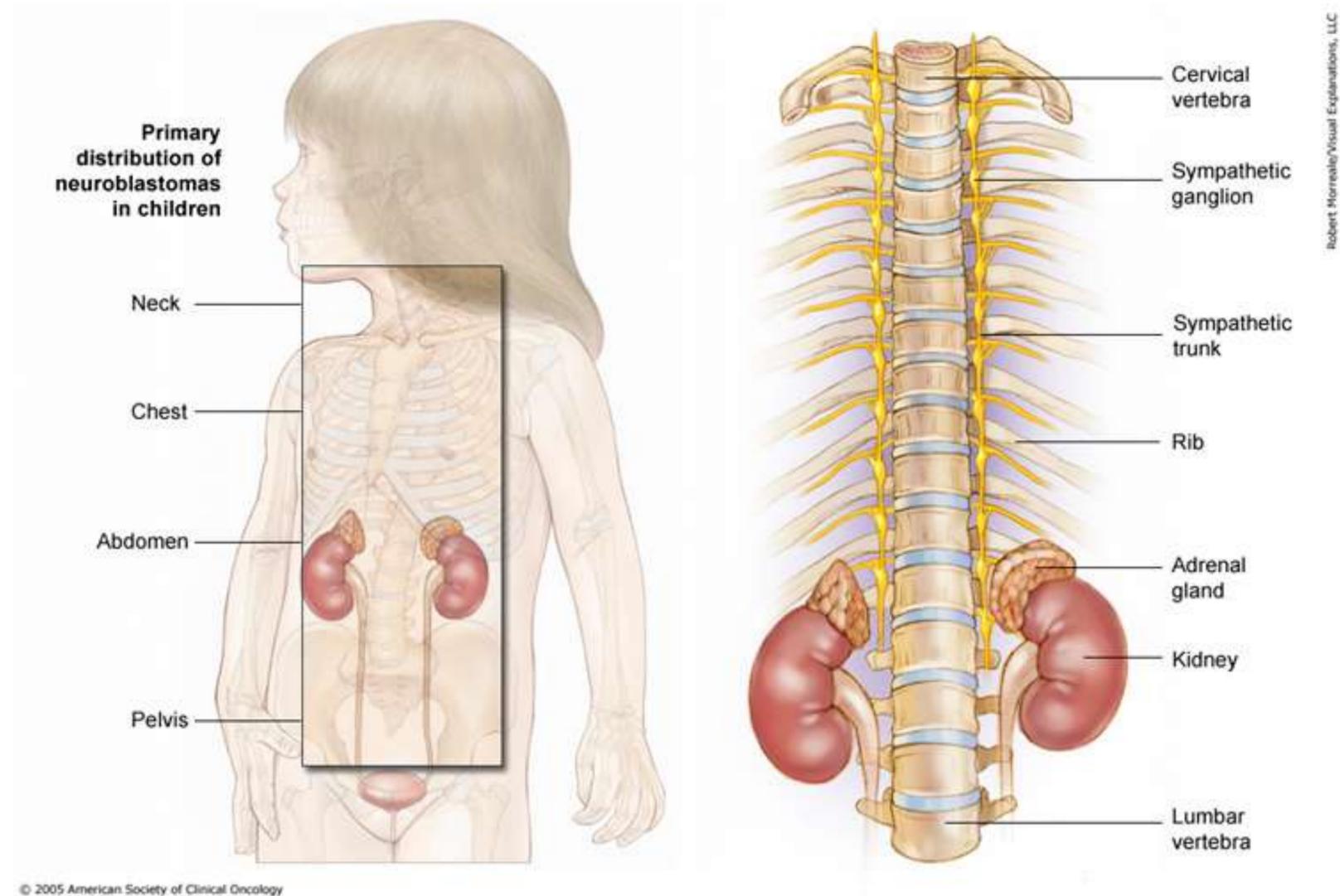
# Children's Oncology Group

1955       Cooperative group system for clinical research

Pediatric Oncology Group (POG)
Children's Cancer Group (CCG)
National Wilms' Tumor Study Group (NWTS)
Intergroup Rhabdomyosarcoma Study Group (IRSG)

2000       Children's Oncology Group (COG)

# Neuroblastoma



Primary distribution of neuroblastomas in children

Neck
Chest
Abdomen
Pelvis

© 2005 American Society of Clinical Oncology

Cervical vertebra
Sympathetic ganglion
Sympathetic trunk
Rib
Adrenal gland
Kidney
Lumbar vertebra

Robert Morreale/Visual Explanations, LLC

The most common solid tumor in children

# Neuroblastoma is rare

| Site | New Cases 2012 |
|---|---|
| Breast | 229,000 |
| Colon | 103,000 |
| Prostate | 241,000 |
| Lung | 226,000 |
| Neuroblastoma | 800 |

Studying such a rare tumor is especially difficult and requires collaborative, consortium-driven research.

# International Neuroblastoma Research Group

Established in 2004



Children's Oncology Group (COG)
4235

Germany
1938

Japan
470

SIOPEN
936

Charged with collecting patients and standardizing terminologies

| Group | Number |
|-------|--------|
| COG | 4235 |
| Germany | 1938 |
| Japan | 470 |
| SIOPEN | 936 |
| Total | 8800 |

The largest collection of neuroblastoma patients and has resulted in more than a dozen high-impact papers

# Neuroblastoma consensus data dictionary

| Field Name | Data Type | Description | Value Constraints |
|---|---|---|---|
| INRG_ID | TEXT | Unique Patient identification number, assigned by the iINRGdb staff after data submission | |
| USI | TEXT | Universal specimen index (COG patients) | |
| AGE | INTEGER | Age (in days) on the date of diagnosis | |
| YEAR | TEXT | Year of diagnosis/enrollment (YYYY) | |
| INIT_TREAT | INTEGER | Initial patient treatment | 0=None (observation)<br>1=Surgery alone<br>2=Conventional-dose chemotherapy (2-8 cycles) plus surgery<br>3=Intensive multi-modality therapy: specific type unknown<br>4=Intensive multi-modality therapy: no stem cell or bone marrow transplant<br>5=Intensive multi-modality therapy: plus stem cell or bone marrow transplant<br>6=Intensive multi-modality therapy: plus stem cell or bone marrow transplant and anti-GD2 antibody<br>9=Unknown |
| INSS_STAGE | INTEGER | INSS stage | 1=Stage 1<br>2=Stage 2a<br>3=Stage 2b<br>4=Stage 3<br>5=Stage 4<br>6=Stage 4s<br>9=Unknown |
| EVANS_STAGE | INTEGER | | 1=Stage I<br>2=Stage II<br>3=Stage III<br>4=Stage IV<br>5=Stage IVs<br>9=Unknown |
| MYCN | INTEGER | MYCN status | 1=Amplified (> 4 times of the reference on chromosome 2q)<br>0=Not amplified (≤ 4 times of the reference on chromosome 2q)<br>9=Unknown, not done, unsatisfactory, in progress |
| PLOIDY | INTEGER | Ploidy | 1=DNA Index ≤ 1 (hypodiploid, diploid)<br>0=DNA Index > 1 (hyperdipolid)<br>9=Unknown |

THE UNIVERSITY OF CHICAGO MEDICINE & BIOLOGICAL SCIENCES | CENTER FOR RESEARCH INFORMATICS

# International Neuroblastoma Research Group

Established in 2004



Children's Oncology Group (COG)
4500

Germany
1900

Japan
500

SIOPEN
900

The good news

8800 patients!

The bad news?

# Neuroblastoma Commons Cohort Discovery

# Demo query 1

Stage 4, 18m-5y, TARGET data

# Neuroblastoma Commons Cohort Discovery



Stage 4, 18m-5y, TARGET data

# Neuroblastoma Commons Cohort Discovery



Stage 4, 18m-5y, TARGET data

# Neuroblastoma Commons Cohort Discovery



Stage 4, 18m-5y, TARGET data

# Neuroblastoma Commons Cohort Discovery



Links to external data sets

# Neuroblastoma Commons Cohort Discovery

# Demo query 2

Favorable biology, tissue available

# Neuroblastoma Commons Cohort Discovery



Favorable biology, tissue available

THE UNIVERSITY OF CHICAGO MEDICINE & BIOLOGICAL SCIENCES | CENTER FOR RESEARCH INFORMATICS

# Neuroblastoma Commons Cohort Discovery



Links to external data sets

# Neuroblastoma Commons Cohort Discovery



Favorable biology, tissue available

# Neuroblastoma Data Commons



INRG cohort discovery

Command-line analysis

Genomic data commons search

Visualization

# Data growth

| Year | COG | SIOPEN | GPOH | Japan | Total |
|------|------|--------|------|-------|-------|
| 2004 | 4235 | 2157 | 1938 | 470 | **8800** |
| 2012 | 6127 | 2504 | 1938 | 470 | **11039** |
| 2013 | 11642 | 2504 | 1938 | 470 | **16554** |
| 2015 | 13060 | 2504 | 1938 | 470 | **17972** |
| 2016 | 13937 | 2664 | 1938 | 470 | **19009** |

# Neuroblastoma data commons

# Federated Authentication

# Governance / Regulation / Compliance

## Connect to Bionimbus Resources

**✔ eRA Commons Account Found!**

SVOLCHEN

**✔ All Required Certificates Found!**

You have already provided all of the required documentation to gain access to Bionimbus resources.

**✔ TARGET Project Access Found!**

TARGET Project Access Found

**✔ Storage Resources Granted!**

Storage resources have been granted to this user.

**✔ Select Bionimbus Storage Bucket**

**Please select the bucket where you'd like to place the data:** svolchen ⬍

Next

THE UNIVERSITY OF CHICAGO MEDICINE & BIOLOGICAL SCIENCES | CENTER FOR RESEARCH INFORMATICS

# Paradigm for building a pediatric cancer commons

1. Engage cooperative group(s)
2. Define scope
3. Identify funding source
4. Identify infrastructure
5. Engage project team
6. Identify data sources
7. Establish governance, create policies and procedures

8. Create contributor / use agreements
9. Create standards working group to create data dictionary, map elements
10. Create database
11. Build front-end query engine
12. Create and execute communication and education plans
13. Create sustainability model

THE UNIVERSITY OF CHICAGO MEDICINE & BIOLOGICAL SCIENCES | CENTER FOR RESEARCH INFORMATICS

# Governance / Regulation / Compliance

**NEUROBLASTOMA DATA CONTRIBUTOR AGREEMENT**

This Neuroblastoma Data Cloud Agreement (this "**Agreement**") is made as of [DATE] (the "**Effective Date**"), by and between The University of Chicago (the "**University**"), and [PARTNER], a [JURISDICTION OF INCORPORATION] [ENTITY TYPE], [ADDRESS] ("**Partner**"), and, together with the University, the "**Parties**").

**RECITALS**

WHEREAS, the University has created a technology platform (the "**Platform**"), including software, hardware, and other technologies, for storing and harmonizing massive data sets of genomic, electronic medical record, and other information related to neuroblastoma ("N

WHEREAS, as part of the Platform, the University owns and operate authorized researchers and other users with access to Neuroblastoma Dat contributors;

WHEREAS, Partner has assembled large data sets of Neuroblas anonymous individuals and associated clinical data ("**Clinical Data**");

WHEREAS, Partner desires to: (i) contribute certain of its Neuroblas (the "**Contributed Data**"), as further described on one or more Contribute below), to the Platform and (ii) permit the University to provide researchers Contributed Data, subject to the restrictions set forth in this Agreement; and

WHEREAS, the University is willing to accept such Contributed Dat

---

**INTERNATIONAL NEUROBLASTOMA RISK GROUP**
**MASTER DATA USE AGREEMENT**

This International Neuroblastoma Risk Group Data Use Agreement (this "**Agreement**") is made as of [DATE] (the "**Effective Date**"), by and between The University of Chicago (the "**University**"), and [PARTNER], a [JURISDICTION OF INCORPORATION] [ENTITY TYPE], [ADDRESS] ("**Partner**"), and, together with the University, the "**Parties**").

**RECITALS**

WHEREAS, the University has created a technology platform (the "**Platform**"), including software, hardware, and other technologies, for storing and harmonizing massive data sets of genomic, electronic medical record, and other information related to neuroblastoma;

WHEREAS, as part of the Platform and in collaboration with the International Neuroblastoma Risk Group ("**INRG**") the University owns and operates a data service that provides authorized researchers and other users with access to such genomic, electronic medial record and other information ("**Contributed Data**") provided by various data contributors (each a "**Data Contributor**");

WHEREAS, Partner desires to permit its researchers to access the Contributed Data, subject to the restrictions set forth in this Agreement; and

WHEREAS, the University is willing to provide such access subject to the terms and conditions set forth in this Agreement.

THE UNIVERSITY OF CHICAGO MEDICINE & BIOLOGICAL SCIENCES | CENTER FOR RESEARCH INFORMATICS

# Building pediatric cancer commons

- Acute myeloid leukemia
- Acute lymphoblastic leukemia
- Bone tumors
- Central nervous system tumors
- Germ cell tumors

- Hodgkin Disease
- Neuroblastoma
- Non-Hodgkin lymphoma
- Renal tumors
- Soft tissue sarcoma

# Building a soft-tissue sarcoma data commons

# Paradigm for building a pediatric cancer commons

1. Engage cooperative group(s)
2. Define scope
3. Identify funding source
4. Identify infrastructure
5. Engage project team
6. Identify data sources
7. Establish governance, create policies and procedures
8. Create contributor / use agreements
9. Create standards working group to create data dictionary, map elements
10. Create database
11. Build front-end query engine
12. Create and execute communication and education plans
13. Create sustainability model

# The importance of standards

Pediatric rhabdomyosarcoma

| | | Current | |
|---|---|---|---|
| | Code | Description | |
| SEX | 1 | Male | |
| | 2 | Female | |
| HISTOLOGY | 1 | Alveolar rhabdomyosarcoma | |
| | 2 | Embryonal rhabdomyosarcoma | |
| | 3 | Botryoid rhabdomyosarcoma | |
| | 4 | Not otherwise specified | |
| | 5 | Undifferentiated sarcoma | |
| | 6 | Sarcoma, not classifiable | |
| | 7 | Spindle cell sarcoma | |
| | 8 | Ectomesenchymoma | |
| | 9 | Other | |
| | 10 | Mixed rhabdomyosarcoma | |
| | 99 | Unknown | |

# The importance of standards

Pediatric rhabdomyosarcoma

| | Current | | Proposed | | |
|---|---|---|---|---|---|
| | Code | Description | Code | Description | Source |
| **SEX** | 1 | Male | C20197 | Male | NCI |
| | 2 | Female | C16576 | Female | |
| | | | C45908 | Intersex | |
| | | | C17998 | Unknown | |
| **HISTOLOGY** | 1 | Alveolar rhabdomyosarcoma | 404053004 | Alveolar rhabdomyosarcoma | SNOMED |
| | 2 | Embryonal rhabdomyosarcoma | 404051002 | Embryonal rhabdomyosarcoma | |
| | 3 | Botryoid rhabdomyosarcoma | 404052009 | Botryoid rhabdomyosarcoma | |
| | 4 | Not otherwise specified | 302847003 | Rhabdomyosarcoma | |
| | 5 | Undifferentiated sarcoma | 128734000 | Undifferentiated sarcoma | |
| | 6 | Sarcoma, not classifiable | 302847003 | Rhabdomyosarcoma | |
| | 7 | Spindle cell sarcoma | 9801004 | Spindle cell sarcoma | |
| | 8 | Ectomesenchymoma | 128750008 | Rhabdomyosarcoma with ganglionic differentiation | |
| | 9 | Other | -- | -- | |
| | 10 | Mixed rhabdomyosarcoma | 62383007 | Mixed type rhabdomyosarcoma | |
| | 99 | Unknown | C17998 | Unknown | NCI |

# Paradigm for pediatric cancer commons

# Gabriella Miller Kids First Pediatric Data Resource Center

# Paradigm for pediatric cancer commons

# Acknowledgements



Center for research informatics    Center for Data-intensive science