# Comparative Methods of CCSG Data Gathering:
# Biosketches, Publications, Grants

**A look at the Past , Present, and Future  of Data Collection**

CCAF-IT 2017

Ben Busby, Mahendra Yatawara, & Susan Sharpe

NIH U.S. National Library of Medicine

CEDARS-SINAI®

MOFFITT CANCER CENTER®

NCBI

# Biosketches:

## sometimes member data collection is like herding cats...

Susan Sharpe, MA

CEDARS-SINAI®  MOFFITT CANCER CENTER  NCBI
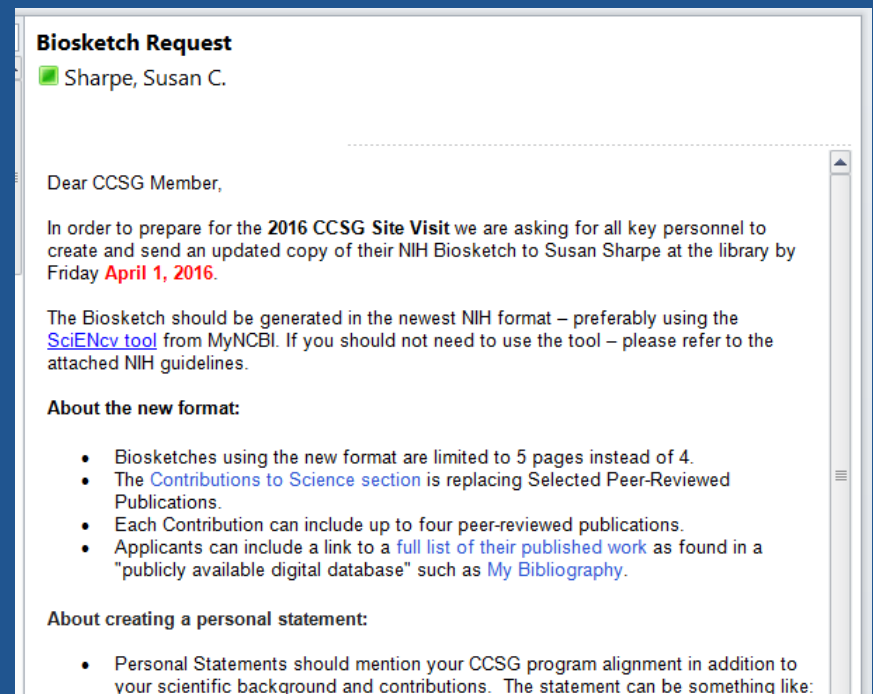
# Biosketches: The expectation

- Routine process, simple 5 page CV of relevant work & interests.

- SciENcv – latest and greatest template.

- Expected Gathering Process: Ask and receive.

**Biosketch Request**

Sharpe, Susan C.

Dear CCSG Member,

In order to prepare for the **2016 CCSG Site Visit** we are asking for all key personnel to create and send an updated copy of their NIH Biosketch to Susan Sharpe at the library by Friday **April 1, 2016**.

The Biosketch should be generated in the newest NIH format – preferably using the SciENcv tool from MyNCBI. If you should not need to use the tool – please refer to the attached NIH guidelines.

**About the new format:**

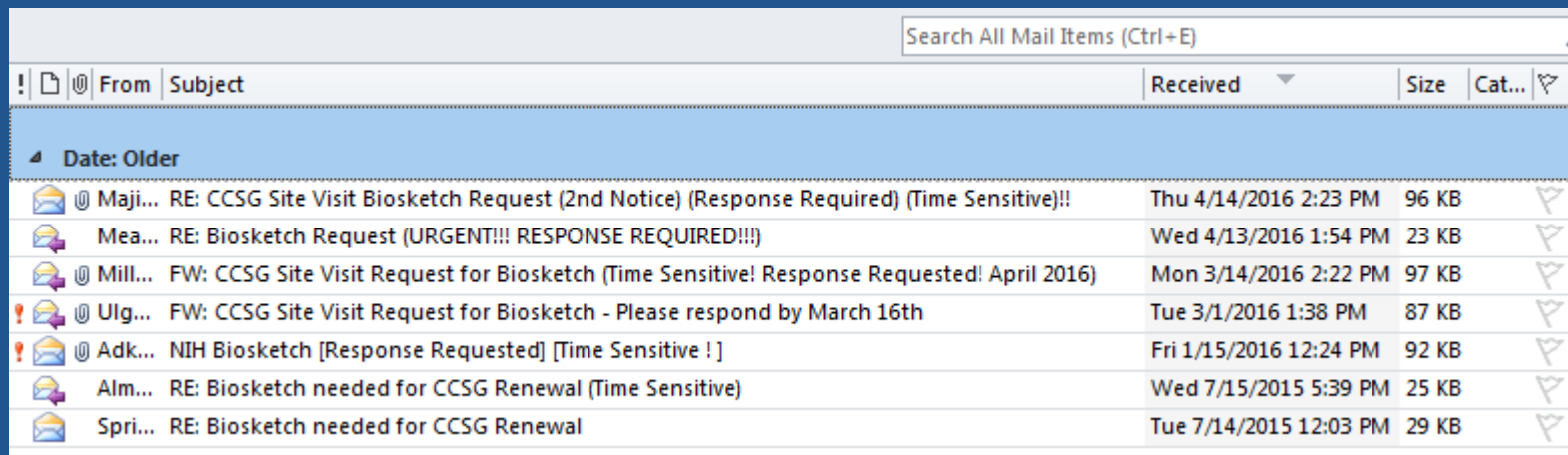- Biosketches using the new format are limited to 5 pages instead of 4.
- The Contributions to Science section is replacing Selected Peer-Reviewed Publications.
- Each Contribution can include up to four peer-reviewed publications.
- Applicants can include a link to a full list of their published work as found in a "publicly available digital database" such as My Bibliography.

**About creating a personal statement:**

- Personal Statements should mention your CCSG program alignment in addition to your scientific background and contributions. The statement can be something like:

U.S. National Library of Medicine

CEDARS-SINAI

MOFFITT CANCER CENTER

NCBI

# The reality:
# a process firmly rooted in the past

1. Email all required personnel request for Biosketch. Include links to SciENcv, provide Word Template, & latest instructions.
2. Wait. Some biosketches return. Edit. Store locally or send to shared drive.
3. Email personnel request for Biosketch Reminder. Add high priority message to email.
4. Wait. Some biosketches return. Edit. Store locally or send to shared drive.
5. Send messages to Faculty Leaders asking for support and encouragement.
6. Wait. Some biosketches return. Edit. Store locally or send to shared drive.
7. Rinse-Repeat x10 times.
8. Biosketches are gathered. Review and make final edits.

| ! | 🗋 | ⑩ | From | Subject | Received | ▼ | Size | Cat... | ⛿ |
|---|---|---|------|---------|----------|---|------|--------|---|
| | | | | ▲ **Date: Older** | | | | | |
| | | ⑩ | Maji... | RE: CCSG Site Visit Biosketch Request (2nd Notice) (Response Required) (Time Sensitive)!! | Thu 4/14/2016 2:23 PM | | 96 KB | | ⛿ |
| | | | Mea... | RE: Biosketch Request (URGENT!!! RESPONSE REQUIRED!!!) | Wed 4/13/2016 1:54 PM | | 23 KB | | ⛿ |
| | | ⑩ | Mill... | FW: CCSG Site Visit Request for Biosketch (Time Sensitive! Response Requested! April 2016) | Mon 3/14/2016 2:22 PM | | 97 KB | | ⛿ |
| ❗ | | ⑩ | Ulg... | FW: CCSG Site Visit Request for Biosketch - Please respond by March 16th | Tue 3/1/2016 1:38 PM | | 87 KB | | ⛿ |
| ❗ | | ⑩ | Adk... | NIH Biosketch [Response Requested] [Time Sensitive ! ] | Fri 1/15/2016 12:24 PM | | 92 KB | | ⛿ |
| | | | Alm... | RE: Biosketch needed for CCSG Renewal (Time Sensitive) | Wed 7/15/2015 5:39 PM | | 25 KB | | ⛿ |
| | | | Spri... | RE: Biosketch needed for CCSG Renewal | Tue 7/14/2015 12:03 PM | | 29 KB | | ⛿ |

Search All Mail Items (Ctrl+E)

**Names redacted to protect the guilty.**

U.S. National Library of Medicine

CEDARS-SINAI®

MOFFITT CANCER CENTER®

NCBI

# SciENcv:

## What went right:

- Automatically puts information in new format
- Create multiple versions
- Share entry and upkeep responsibilities with delegates
- Create sharable URL
- Links to MyBibliography

## What went wrong:

- URL version doesn't enable viewers to download.
- De-centralized management (PI-centric, instead of institutionally)
- No delivery mechanism:
  - PDFs & Emails can be lost, forgotten, deleted, etc.

NIH U.S. National Library of Medicine

NCBI

# How do members of CCAF gather Biosketches?

- 69% of respondents rely on Members to submit and maintain Biosketches

- 14% have homegrown systems that centralize and keep track of Biosketches

- 4% have some sort of vendor system

| Biosketch Methods of Collection: | # of Responses: |
|---|---|
| We rely on **Members** to submit and maintain Biosketches. | 34 |
| We use a **homegrown system** to collect, create, manage, & store. | 7 |
| **Other:** Members write, we edit or provide templates. | 4 |
| **Homegrown Other:** Yes, we have a homegrown solution, but… | 1 |
| **Vendor Other:** Yes, we have a vendor solution, but… | 1 |
| We use a **vendor supported** system to collect, create, manage, & store. | 1 |
| We use existing **NIH** provided tools (NCBI, etc). | 1 |
| **Grand Total** | **49** |

# Are we happy?

| Satisfaction | # of Responses: |
|---|---|
| Dissatisfied | 16 |
| OK | 16 |
| Satisfied | 12 |
| Very Dissatisfied | 2 |
| Very Satisfied | 3 |
| Grand Total | 49 |

- 63% are pretty OK with current methods
- 37% are not

| Who's Happy Here: | #of Responses |
|---|---|
| **OK** | **16** |
|   Other: Members write, we edit. | 3 |
|   Vendor Other: Yes, we have a vendor solution, but… | 1 |
|   We rely on Members to submit and maintain Biosketches. | 12 |
| **Satisfied** | **12** |
|   We rely on Members to submit and maintain Biosketches. | 7 |
|   We use a homegrown system to collect, create, manage, & store. | 4 |
|   We use existing NIH provided tools (NCBI, etc). | 1 |
| **Very Satisfied** | **3** |
|   We use a homegrown system to collect, create, manage, & store. | 3 |
| Grand Total | 31 |



PAIN SCALE

# What's working?

## Vendor Products

- Complion

## Other:

- Centralized department devoted to entry & management

- Regular (Monthly!) updates

## Homegrown Products

- Mission-Based Management

- Nexus

- Faculty Collaboration Database (FCD)

- Customized SciENcv Clone

# The times they aren't changing....

| Planning on Changing Anytime Soon? | # of Responses: |
|---|---|
| No. | 36 |
| We want to change, but have no plans. | 1 |
| We're curious about what others do. | 1 |
| We're looking towards our University to implement a process/product. | 2 |
| Yes. Attempting to choose between homegrown & vendor solutions. | 1 |
| Yes. We have plans to move to vendor solution. | 4 |
| Yes. We are buying a vendor solution. | 1 |
| Yes. We are working on a homegrown solution. | 2 |
| Grand Total | 48 |

- 16% are planning on changing their methods
- 4% waiting for the next big thing

**Vendors:**

- Research Management System (RMS)
- Nexus
- Salesforce
- Mendix
- Café
- RES Forte

# Publications

**How Moffitt collects Pub Data:**

- Nightly search of author names via API to MEDLINE
- Download into holding queue
  - Impact Factor automatically assigned
- Daily author verification screening by human



| | Pending | Add | Reports | Weekly Report | Library Reports | Shared Resources | Staff Tagging Em... |
|---|---|---|---|---|---|---|---|

**1555 Journals Pending Import**

| View | Year | Month | PMID | Citation |
|---|---|---|---|---|
| VIEW | 2008 | Dec | 18682882 | Sanchez JA, Vogel JD, Kalady MF, Bronner MP, |
| VIEW | 2008 | Dec | 18930709 | Kado M, Lee JK, Hidaka K, Miwa K, Murohara T Dec;377(2):413-418. Pubmedid: 18930709. |

# What works for us, may not work for you:

**Pros:**

- Automated & customizable search algorithm
- Very little need for author input
- Standardized citation information
- Ability to pull corresponding data: Grant IDs, ORCiD, MeSH, IF

**Cons:**

- Labor intensive
- Centralizing Screening process requires dedicated staff members
- Author Name Disambiguation remains a stumbling block

# Publications: The current state and a look at our center's process

Mahendra Yatawara, MBA

# CCAF-IT 2017 Survey

- http://moffitt.libsurveys.com/CCAFData
- Survey sent out April 20[th]
- **Survey active until May 2[nd]**
- Institutions responding: 44

# Institutions Responding

# How do Centers manage Pubs for CCSG?

## Publications Systems



Pie chart legend:
- Vendor System
- Home Grown System
- Other
- NIH Tools
- Managed by Members

Chart values:
- 0, 0%
- 7, 14%
- 7, 14%
- 12, 25%
- 23, 47%

## Vendor and Other

| CAFÉ by USC | 3 |
|---|---|
| Opus/EVAL by Forte | 2 |
| Lattice Grid | 2 |

| Homegrown & NIH |
|---|
| Nexus |
| Homegrown |

# Satisfaction with current Pubs System?

## Publications Systems

0, 0%

7, 14%    7, 14%

12, 25%

23, 47%

- ☐ Vendor System
- ☐ Home Grown System
- ☐ Other
- ☐ NIH Tools
- ☐ Managed by Members

## Satisfaction Level

3, 6%

8, 15%

12, 23%

14, 27%

15, 29%

- ☐ Very Satisfied
- ☐ Satisfied
- ☐ OK
- ☐ Dissatisfied
- ☐ Very Dissatisfied

# Plans to change Pubs Solution?

## Satisfaction and Change



## Possible alternatives

# How Cedars-Sinai Collects Publications Data

Initial pubs import

| Public Bibliography CREATED in NCBI Portal [Member] | → | Public bibliography link SAVED in CMAPS [Member / Admin staff] | → | Pubs imported via single click in CMAPS [Member / Admin staff] | → | Impact factors auto-tagged to publication upon import [System] | → | Cores utilized and Cancer relevance assigned to publications [Member] |

# How Cedars-Sinai Collects Publications Data

**Initial pubs import**

Public Bibliography CREATED in NCBI Portal [Member] → Public bibliography link SAVED in CMAPS [Member / Admin staff] → Pubs imported via single click in CMAPS [Member / Admin staff] → Impact factors auto-tagged to publication upon import [System] → Cores utilized and Cancer relevance assigned to publications [Member]

**Subsequent pubs import**

Public Bibliography UPDATED in NCBI Portal [Member] → Pubs imported via single click in CMAPS [Member / Admin staff] → Impact factors auto-tagged to publication upon import [System] → Cores utilized and Cancer relevance assigned to publications [Member]

# How is this process working for us?

## Key Advantages

- Reduction in non-value added work from CC Admin
- Members maintain in single location (NCBI portal)
- Auto-assignment of Impact Factor
- One-click reports

## Limitations

- Reminders for Members to keep NCBI Bibliography up-to-date
- Reminders for Members to allocate Core usage and Cancer Relevance to pubs in CMAPS

# The Futures: Biosketches, Grants, Pubs... and Data!

Ben Busby, NCBI

# NCBI

Saved Searches

My Bibliography

Collections

SciENcv

## Search NCBI databases

Search : PubMed ▾

[_____]  Search

Hint: clicking the "Search" button without any terms listed in the search box will transport you to that database's homepage.

## My Bibliography

Your bibliography contains 1 items.

Share your bibliography with this URL:
http://www.ncbi.nlm.nih.gov/sites/myncbi/1Rgszs64zgoAt/bibliography/46426933/public/?sort=date&direction=descending

**Most recent citations:**

Bocik. TEST!. Testing Journal. 2015;

Manage My Bibliography »

## Recent Activity

| Time | Database | Type | Term |
|------|----------|------|------|
| 8:51 PM | Books | record | My Bibliography - My NCBI Help |
| 8:49 PM | Books | record | SciENcv - My NCBI Help |
| 27-Feb-2015 | Assembly | search | txid9887[Organism] |
| 27-Feb-2015 | Nucleotide | record | Muntiacus muntjak vaginalis clone C... |
| 27-Feb-2015 | Nucleotide | record | Muntiacus muntjak vaginalis clone I... |
| 27-Feb-2015 | Nucleotide | search | txid9887[Organism] AND (biomol_geno... |
| 27-Feb-2015 | BioSample | record | Indian muntjac whole genome BAC lib... |
| 27-Feb-2015 | BioSample | search | Muntjak |

## Saved Searches

You don't have any saved searches yet.

Go and create some saved searches in PubMed or our other databases.

Manage Saved Searches »

## Collections

| Collection Name | | Items | Settings/Sharing | Type |
|-----------------|------|-------|------------------|------|
| Favorites | edit | 0 | ⚙ Private | Standard |
| My Bibliography | edit | 1 | ⚙ Public | Standard |
| Other Citations | edit | 0 | ⚙ Private | Standard |

Manage Collections »

## Filters

Filters for: PubMed ▾

You do not have any active filters for this database.
Add filters for the selected database.

Manage Filters »

## SciENcv

| Name | Last Update | Sharing | Type |
|------|-------------|---------|------|
| NewSketch No External | 09-Feb-2015 | Private | NIH Biosketch |
| ORCIDTEST | 09-Feb-2015 | Private | Old NIH Biosketch |
| Test2 | 06-Mar-2015 | Private | Old NIH Biosketch |

Manage SciENcv »

# NCBI

# Better PubMed Searches!

# For more information go to:
## ncbi.nlm.nih.gov/learn

# E-Utilities (Eutils)

| Entrez Database | UID common name | E-utility Database Name |
|---|---|---|
| BioProject | BioProject ID | bioproject |
| BioSample | BioSample ID | biosample |
| Biosystems | BSID | biosystems |
| Books | Book ID | books |
| Conserved Domains | PSSM-ID | cdd |
| dbGaP | dbGaP ID | gap |
| dbVar | dbVar ID | dbvar |
| Epigenomics | Epigenomics ID | epigenomics |
| EST | GI number | nucest |
| Gene | Gene ID | gene |
| Genome | Genome ID | genome |
| GEO Datasets | GDS ID | gds |
| GEO Profiles | GEO ID | geoprofiles |
| GSS | GI number | nucgss |
| HomoloGene | HomoloGene ID | homologene |
| MeSH | MeSH ID | mesh |
| NCBI C++ Toolkit | Toolkit ID | toolkit |
| NCBI Web Site | Web Site ID | ncbisearch |
| NLM Catalog | NLM Catalog ID | nlmcatalog |
| Nucleotide | GI number | nuccore |

| | | |
|---|---|---|
| PopSet | PopSet ID | popset |
| Probe | Probe ID | probe |
| Protein | GI number | protein |
| Protein Clusters | Protein Cluster ID | proteinclusters |
| PubChem BioAssay | AID | pcassay |
| PubChem Compound | CID | pccompound |
| PubChem Substance | SID | pcsubstance |
| PubMed | PMID | pubmed |
| PubMed Central | PMCID | pmc |
| SNP | rs number | snp |
| SRA | SRA ID | sra |
| Structure | MMDB-ID | structure |
| Taxonomy | TaxID | taxonomy |
| UniGene | UniGene Cluster ID | unigene |

# *Introducing…* Entrez Direct
# The E-utilities on the UNIX command line

```
esearch -db gene -query "foxp2[gene]
AND human[orgn]" | \

elink -target protein -name
gene_protein_refseq | \

efetch -format fasta
```

ftp.ncbi.nlm.nih.gov/entrez/entrezdirect/

NCBI

# *The EDirect Cookbook!*

**Convert article DOI to PMID**

Description (optional):
Written by: NCBI Folks (12/14/2016)
Confirmed by: Mike Davidson (NLM) (12/16/2016, v5.80)
Databases: pubmed

```
esearch -db pubmed -query "10.1111/j.1468-3083.2012.04708.x" | \
esummary | \
xtract -pattern DocumentSummary -block ArticleId -sep "\t" -tab "\n" -element IdType,Value | \
grep -E '^pubmed|doi'
```

**Access organism specific meta-data from NCBI genome database**

Description (optional):
Written by: NCBI Folks (12/14/2016)
Confirmed by:
Databases: genome, bioproject

```
esearch -db genome -query "22954[uid]" | \
elink -target bioproject | \
efetch -format xml | \
xtract -pattern DocumentSummary -element Salinity OxygenReq OptimumTemperature TemperatureRange Habitat
```

**Get the status of records from PubMed search**

Description (optional):
Written by: NCBI Folks (12/14/2016)
Confirmed by: Mike Davidson (NLM) (12/16/2016, v5.80)
Databases: pubmed

```
esearch -db pubmed -query "pde3a AND 2016[dp]" | \
esummary | \
xtract -pattern DocumentSummary -element Id RecordStatus
```

**Conduct a PubMed search and retrieve the results as a list of PMIDs**

Description (optional):
Written by: Mike Davidson (2/22/2017)
Confirmed by: Mike Davidson (NLM) (2/22/2017, v6.30)
Databases: pubmed

```
esearch -db pubmed -query "seasonal affective disorder" | efetch -format uid
```

**Sort the hits by sequence length in nucleotide database**

Google for
EDirect Cookbook

# BioProject



NCBI    Resources ⌄    How To ⌄

**BioProject**

[ BioProject ⌄ ]    nutrition

Create alert    Advanced

**Project Types**
Umbrella (40)
Primary submission (694)
RefSeq (12)

**Data Types**
Epigenomics (31)
Genome sequencing (38)
Metagenome (28)
Metagenomic assembly (1)
Other (42)
Phenotype/genotype (7)
Random survey (1)
Targeted locus (8)
Transcriptome (474)

**Project Data**
Nucleotide (50)
Protein (32)
Assembly (44)
SRA (187)
GEO DataSets (507)

**Scope**
Monoisolate (112)
Multi-isolate (505)
Multi-species (21)
Environmental (58)
Other (10)

Display Settings: ⌄    Summary, 20 per page, Sorted by Default order          Send to: ⌄

**Search results**

Items: 1 to 20 of 746                 << First  < Prev  Page [1] of 38  Next >  Last >>

☐  **Bacteria**
1.   Bacteria sequenced from reef-building corals Raw sequence reads
     Project data type: Raw sequence reads
     Scope: Multispecies
     University of Hawaii at Manoa
     Accession: PRJNA355371   ID: 355371

☐  **A Novel Regulatory Region for Amylose Synthesis in Rice Grains Identified by Systems Genetics**
2.   **Approach.**
     Organism: Oryza sativa Indica Group
     Taxonomy: *Oryza sativa Indica Group (long-grained rice)*
     Project data type: Transcriptome or Gene expression
     Scope: Multiisolate
     IRRI
     Accession: PRJNA355111   ID: 355111

☐  **panda gut metagenome**
3.   Panda gut fungal metagenome: raw sequence reads
     Taxonomy: *gut metagenome*
     Project data type: Raw sequence reads
     Scope: Environment

NIH⟩ U.S. National Library of Medicine                                    NCBI

BioProject    [BioProject ▼]    nutrition

Create alert    Advanced

Project Types
Umbrella (40)
Primary submission (694)
RefSeq (12)

Data Types
Epigenomics (31)
Genome sequencing (38)
Metagenome (28)
Metagenomic assembly (1)
Other (42)
Phenotype/genotype (7)
Random survey (1)
Targeted locus (8)
Transcriptome (474)

Project Data
Nucleotide (50)
Protein (32)
Assembly (44)
SRA (187)
GEO DataSets (507)

Scope
Monoisolate (112)
Multi-isolate (505)
Multi-species (21)
Environmental (58)
Other (10)

Display Settings    ault order    Send to: ▼

Search resu

Items: 1 to 20    << First    < Prev    Page [1]    of 38    Next >    Last >>

Data Types
Epigenomics (31)
Genome sequencing (38)
Metagenome (28)
Metagenomic assembly (1)
Other (42)
Phenotype/genotype (7)
Random survey (1)
Targeted locus (8)
Transcriptome (474)

Project Data
Nucleotide (50)
Protein (32)
Assembly (44)
SRA (187)
GEO DataSets (507)

Scope
Monoisolate (112)
Multi-isolate (505)
Multi-species (21)
Environmental (58)
Other (10)

☐ Bacteria
1. Bacteria s      sequence reads
   Project data
   Scope: Mult
   University o
   Accession:

☐ A Novel R      is in Rice Grains Identified by Systems Genetics
2. Approach.
   Organism: G
   Taxonomy: G
   Project data
   Scope: Mult
   IRRI
   Accession:

☐ panda gut
3. Panda gut      ds
   Taxonomy: g
   Project data
   Scope: Envi

U.S. National Library of Medicine
NIH
NCBI

# Reporting

# dbGaP

# Minimizing Data Transfer



sam-dump.2.6.3 --aligned-region 17:41243452-41277500
SRR925743 > BRCA1.sam

# Minimizing Data Transfer

# Minimizing Data Transfer

```
[ec2-user@ip-172-16-243-238 Reference]$ hisat2-build `ls *fasta | awk '{printf("%s,",$1)}' | sed -e 's/,$//'` HT2_ID
X
Settings:
  Output files: "HT2_IDX.*.ht2"
  Line rate: 6 (line is 64 bytes)
  Lines per side: 1 (side is 64 bytes)
  Offset rate: 4 (one in 16)
  FTable chars: 10
  Strings: unpacked
```

```
[ec2-user@ip-172-16-243-238 Reference]$ ls -ltr
total 21900
-rw-rw-r-- 1 ec2-user ec2-user 4102082 Mar 21 20:55 LK936442.1.fasta
-rw-rw-r-- 1 ec2-user ec2-user 3161919 Mar 21 20:57 LK936443.1.fasta
-rw-rw-r-- 1 ec2-user ec2-user 1790417 Mar 21 21:03 HT2_IDX.4.ht2
-rw-rw-r-- 1 ec2-user ec2-user      26 Mar 21 21:03 HT2_IDX.3.ht2
-rw-rw-r-- 1 ec2-user ec2-user       8 Mar 21 21:03 HT2_IDX.8.ht2
-rw-rw-r-- 1 ec2-user ec2-user       8 Mar 21 21:03 HT2_IDX.7.ht2
-rw-rw-r-- 1 ec2-user ec2-user 1790424 Mar 21 21:03 HT2_IDX.2.ht2
-rw-rw-r-- 1 ec2-user ec2-user 6581787 Mar 21 21:03 HT2_IDX.1.ht2
-rw-rw-r-- 1 ec2-user ec2-user 1822930 Mar 21 21:03 HT2_IDX.6.ht2
-rw-rw-r-- 1 ec2-user ec2-user 3149021 Mar 21 21:03 HT2_IDX.5.ht2
[ec2-user@ip-172-16-243-238 Reference]$
```

```
[ec2-user@ip-172-16-243-238 SRR3145392]$ hisat2 -f -x ../HT2_IDX --sra  SRR3145392 --no-spliced-alignment --threads
8 > SRR3145392.sam
2868271 reads; of these:
  2868271 (100.00%) were paired; of these:
    465062 (16.21%) aligned concordantly 0 times
    1757830 (61.29%) aligned concordantly exactly 1 time
    645379 (22.50%) aligned concordantly >1 times
    ----
    465062 pairs aligned concordantly 0 times; of these:
      13148 (2.83%) aligned discordantly 1 time
    ----
    451914 pairs aligned 0 times concordantly or discordantly; of these:
      903828 mates make up the pairs; of these:
        737036 (81.55%) aligned 0 times
        130993 (14.49%) aligned exactly 1 time
        35799 (3.96%) aligned >1 times
87.15% overall alignment rate
[ec2-user@ip-172-16-243-238 SRR3145392]$
```

# Minimizing Data Transfer

# Minimizing Data Transfer

# Reporting

# NCBI

# For more information go to:
## ncbi.nlm.nih.gov/learn