

Data Entry: Bane and Boon

Morgan Taschuk

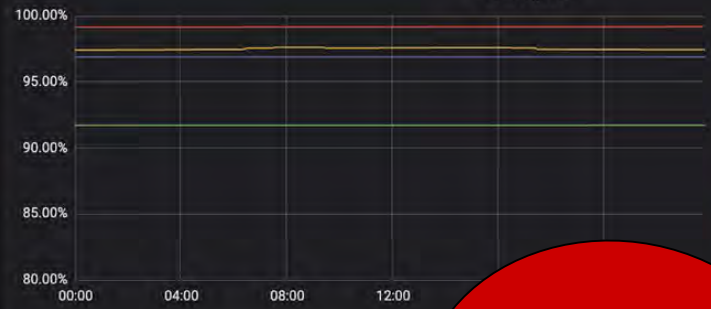
@morgantaschuk



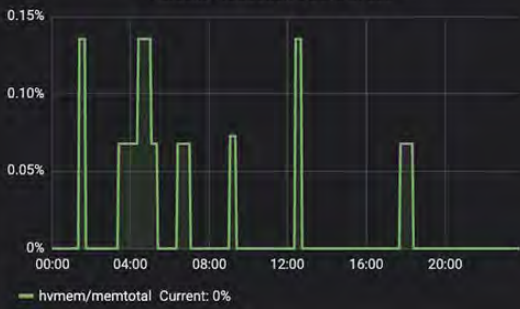
Prod spa...

110 TiB

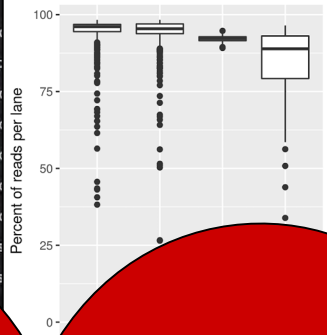
Disk Space



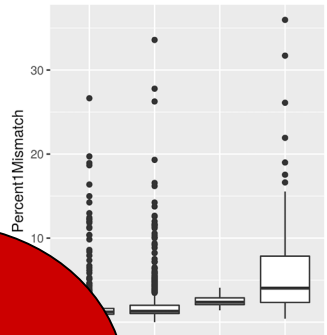
Percent Production Queue in use



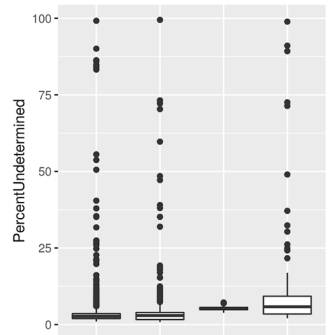
A) 0 mismatches



B) 1 mismatch

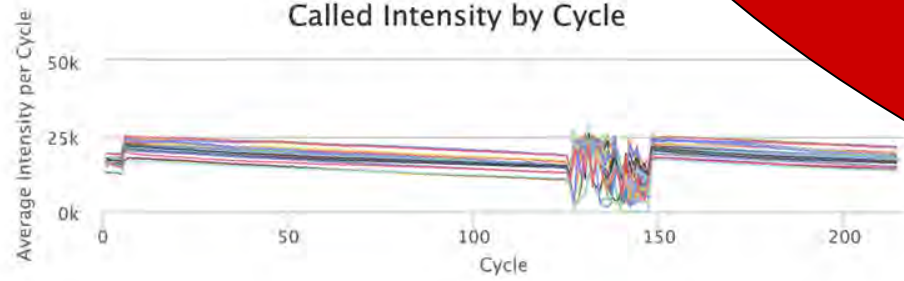


C) Undetermined indices

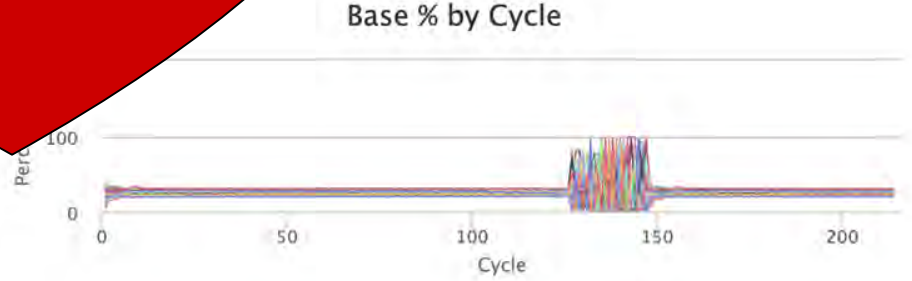


DATA

Called Intensity by Cycle



Base % by Cycle



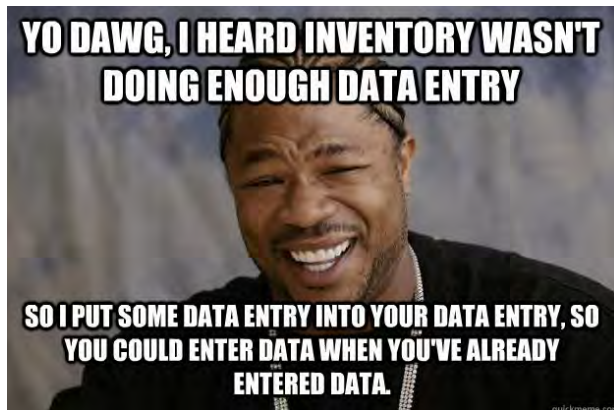
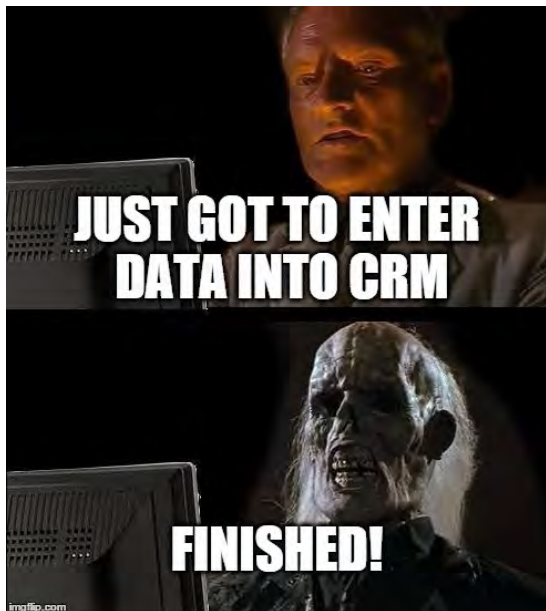
A (Combined) C (Combined) G (Combined) T (Combined)

A (Combined) C (Combined) G (Combined) T (Combined)

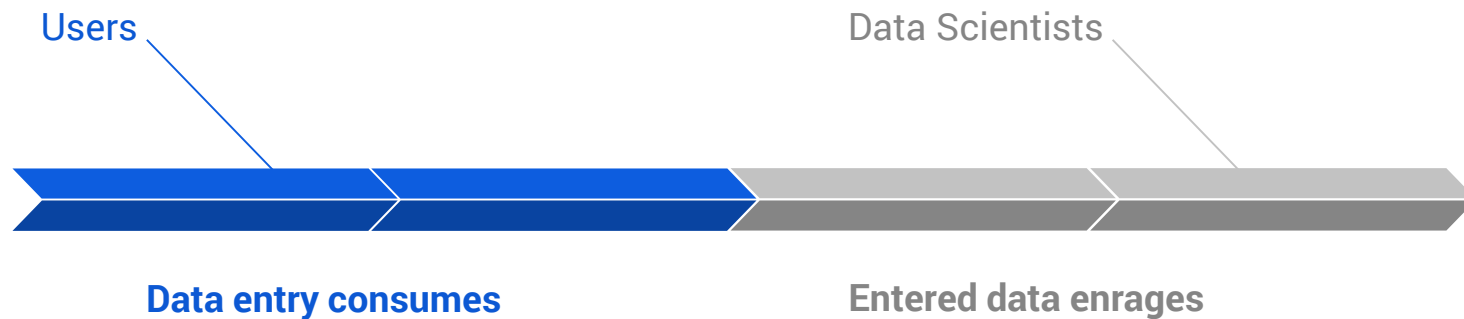
Data entry
drives our
world



Data Entry



Continuum of data entry misery

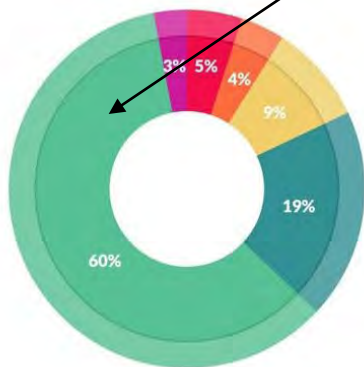




“people in the medical profession actively, viscerally, volubly hate their computers.”

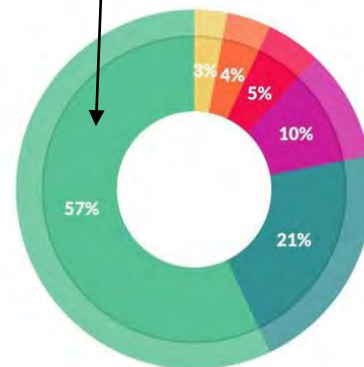
<https://www.newyorker.com/magazine/2018/11/12/why-doctors-hate-their-computers>

Data scientists spend 60% of their time cleaning data and kind of hate it



What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets: 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%



What's the least enjoyable part of data science?

- Building training sets: 10%
- Cleaning and organizing data: 57%
- Collecting data sets: 21%
- Mining data for patterns: 3%
- Refining algorithms: 4%
- Other: 5%

<https://www.forbes.com/sites/gilpress/2016/03/23/data-preparation-most-time-consuming-least-enjoyable-data-science-task-survey-says/>

OICR's data entry continuum

Genomics technicians



Cyanide and Happiness

<https://www.youtube.com/watch?v=zTsGcNpr-k0>

GeneSifter[®] Lab Edition OICR PRODUCTION User: miaschuk Aug 12, 2015 17:04 Get FinchTV Quick Search Go Logout

Home Orders Lab Data System

Template Tracker

- Add Template
- Add Templates (Bulk)
- View Templates
- View Template Worksets
- View Template Reports

Primer Tracker

- Add Primer
- Add Primers (Bulk)
- View Primers

Inventory Tracker

- Add Inventory Item
- Inventory Catalog

Flows

- Flow
- Matrix
- PacE
- Picot
- Slide
- StorE

ContE

- View
- Pool

Instrument Runs

- Add Instrument Run
- View Runs
- View Run Data

Template Details

Edit Template Upload Attachment Set to Archived Derive Children

Go to Workflow Step

Previous: Library Details Current: Ready

OR Ready (current) Go

Enter Next Workflow

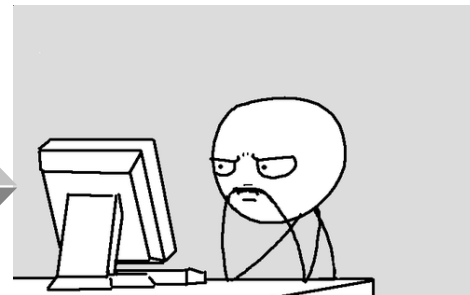
Illumina PE Library Seq Go

Name	CPCG_0154_Pr_P_PE_617_WG	Template Workset	n/a
Template ID	26347	Parents	CPCG_0154_Pr_P_m_1-1_0-2
Template Type	Illumina PE Library	Volume (µL)	n/a
Barcode	0134726728	Concentration (ng/µL)	n/a
Archived	No	Storage Location	n/a

Library Details: Tissue Type	P	Library Details: Library Selection Process	PCR
Library Details: Library Source	GENOMIC	Library Details: Source Template Type	WG
Barcode: Barcode	None		

Laboratory Information Management System (LIMS)

My Team



Users

Data Scientists

Data entry consumes

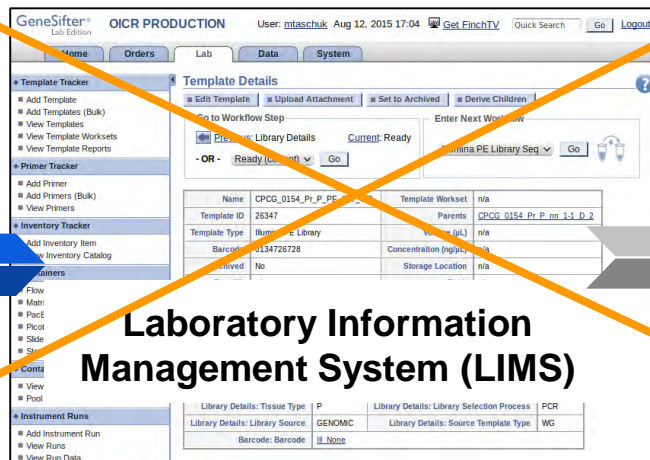
Entered data enrages

Seeing the light in 2015

Genomics

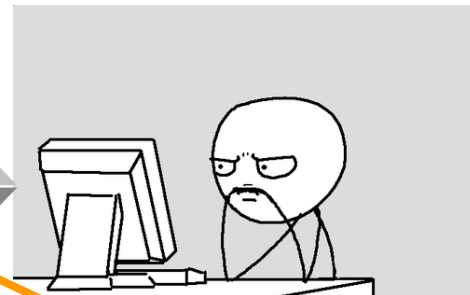


<https://www.youtube.com/watch?v=zTsGcNpr-k0>



Laboratory Information Management System (LIMS)

My Team



Users

Data Scientists

Data entry consumes

Entered data enrages

Meanwhile, in 2015



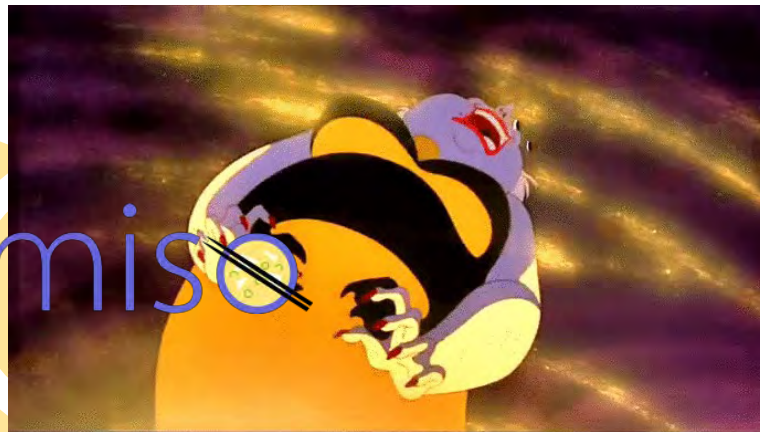
Genomics

Also my team

My Team



Meanwhile, in 2015



Genomics

Also my team

My Team



Let's not repeat our mistakes

Why was the original software so terrible?

GeneSifter[®] Lab Edition OICR PRODUCTION User: mtschuk Aug 12, 2015 17:04 Get FinchTV Quick Search Go Logout

Home Orders Lab Data System

Template Tracker

- Add Template
- Add Templates (Bulk)
- View Templates
- View Template Worksets
- View Template Reports

Primer Tracker

- Add Primer
- Add Primers (Bulk)
- View Primers

Inventory Tracker

- Add Inventory Item
- View Inventory Catalog

Containers

- Flowcell
- Matrix Box
- PacBio
- Picotter Plate
- Slide
- Storage Box

Container Tracker

- View Containers
- Pool Containers

Instrument Runs

- Add Instrument Run
- View Runs
- View Run Data

Template Details

Edit Template Upload Attachment Set to Archived Derive Children

Go to Workflow Step Enter Next Workflow

Previous: Library Details Current: Ready

OR Ready (current) Go Illumina PE Library Seq Go

Name	CPCG_0154_Pr_P_PE_617_WG	Template Workset	n/a
Template ID	26347	Parents	CPCG_0154_Pr_P_nn_1-1_D_1
Template Type	Illumina PE Library	Volume (µL)	n/a
Barcode	0134726728	Concentration (ng/µL)	n/a
Archived	No	Storage Location	n/a
Prep Kit	n/a	Rack	n/a
Description	Created from template CPCG_0154_Pr_P_nn_1_D_2		

Related Links: Details, Workflow History, Results, Relationships, Hierarchy, Attachments

Workflow Data

Library Details: Tissue Origin	Pr	Library Details: Organism	Home
Library Details: Library Type	PE	Library Details: Library Strategy	WGS
Library Details: Tissue Type	P	Library Details: Library Selection Process	PCR
Library Details: Library Source	GENOMIC	Library Details: Source Template Type	WG
Barcode: Barcode	Ill_None		

miso managing information for sequencing operations PRODUCTION

Home My Account My Projects Reports Help

Logged in as: mtschuk Logout

Samples

Show 25 entries Search: Add Sample

Select All

Sample Name	Alias	Sample Class	Type	QC Passed	Location	Last Updated
<input type="checkbox"/> SAM24	AMLXP_0008_Bm_P_nn_1-1_D_51	gDNA (stock)	GENOMIC	Unknown	Unknown	2016-10-07
<input type="checkbox"/> SAM23	AMLXP_0008_Bm_P_nn_1-1	Primary Tumor Tissue	GENOMIC	Unknown	Unknown	2016-10-07
<input type="checkbox"/> SAM22	AMLXP_0008	Identity	GENOMIC	Unknown	Unknown	2016-10-07
<input type="checkbox"/> SAM21	AMLXP_0007_Bm_P_nn_1-1_D_51	gDNA (stock)	GENOMIC	Unknown	Unknown	2016-10-07
<input type="checkbox"/> SAM18	AMLXP_0006_Bm_P_nn_1-1_D_51	gDNA (stock)	GENOMIC	Unknown	Unknown	2016-10-07
<input type="checkbox"/> SAM15	AMLXP_0005_Bm_P_nn_1-1_D_51	gDNA (stock)	GENOMIC	Unknown	Unknown	2016-10-07
<input type="checkbox"/> SAM20	AMLXP_0007_Bm_P_nn_1-1	Primary Tumor Tissue	GENOMIC	Unknown	Unknown	2016-10-07
<input type="checkbox"/> SAM17	AMLXP_0006_Bm_P_nn_1-1	Primary Tumor Tissue	GENOMIC	Unknown	Unknown	2016-10-07
<input type="checkbox"/> SAM14	AMLXP_0005_Bm_P_nn_1-1	Primary Tumor Tissue	GENOMIC	Unknown	Unknown	2016-10-07
<input type="checkbox"/> SAM19	AMLXP_0007	Identity	GENOMIC	Unknown	Unknown	2016-10-07

Tracking

- Samples
- Libraries
- Pools
- Orders
- Sequencing Containers
- Boxes
- Plates
- Sequencers
- Kits
- Indices
- Experiments
- Studies

User Administration

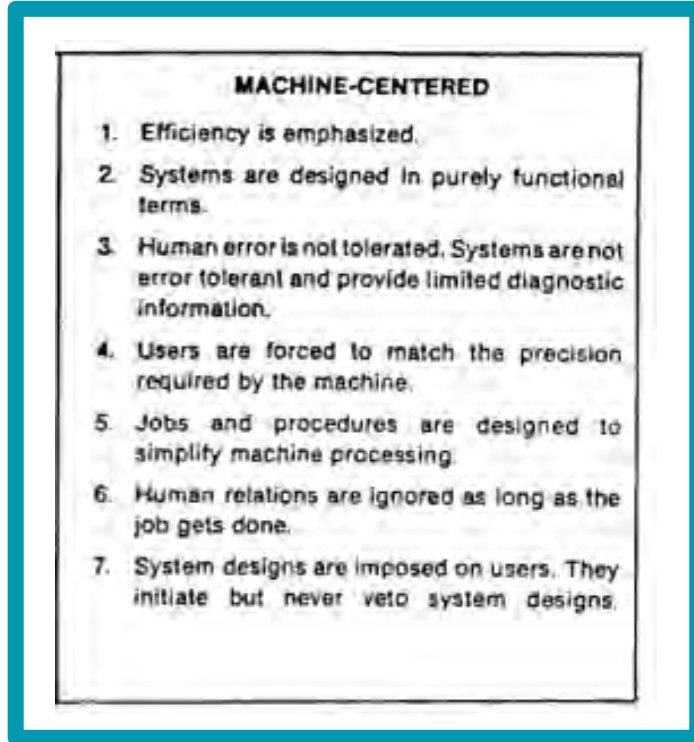
- Users
- Groups

Why does software fail?

“Unfortunately, many computer programs that are well designed in terms of technical criteria, such as run-time efficiency, fail to meet the human or organizational needs they were expected to serve.”

Rob Kling, Towards a person-centered computer technology, 1973

doi: 10.1145/800192.805740



Rob Kling, The Organizational Context of User-Centred Software Designs, 1977
doi: 10.2307/249021

Humans are fallible



MACHINE-CENTERED

1. Efficiency is emphasized.
2. Systems are designed in purely functional terms.
3. Human error is not tolerated. Systems are not error tolerant and provide limited diagnostic information.
4. Users are forced to match the precision required by the machine.
5. Jobs and procedures are designed to simplify machine processing.
6. Human relations are ignored as long as the job gets done.
7. System designs are imposed on users. They initiate but never veto system designs.

USER-CENTERED

1. Systems are valued that increase personal competence and pride in work.
2. People are accepted as non-rational and error-prone.
3. Jobs are designed to be personally satisfying. Automated procedures are designed to fit job needs.
4. The burden of precision is placed on the machine. Systems are forgiving.
5. Users easily obtain/create systems that meet their needs.
6. Users can initiate, veto, and collaborate in system designs.
7. Designs and assumptions are intelligible to users through appropriate technique (modular structures) and clear documentation.

Rob Kling, The Organizational Context of User-Centred Software Designs, 1977
doi: 10.2307/249021

Autocorrect



How do we help our users succeed?



PRODUCTION

[Help](#) | [Report a problem](#) | Logged in as: [mtaschuk](#) | [Logout](#)

Preparation

[Projects](#)
[Samples](#)
[Libraries](#)
[Dilutions](#)
[Worksets](#)
[Pools](#)
[Boxes](#)

Instrument Runs

[Orders](#)
[All](#)
[Active](#)
[Pending](#)
[Sequencing](#)
[Containers](#)
[Runs](#)
[Array Scanning](#)
[Arrays](#)
[Runs](#)
[Instruments](#)

Tools

[Index Distance](#)

Samples

[Create](#) [Edit](#) [Propagate](#) [Print Barcode\(s\)](#) [Download](#) [Parents](#) [Children](#) [Add QCs](#) [Edit QCs](#) [Add to Workset](#) [Attach Files](#)

Show 25 entries

 Search:

	Name	Alias	Sample Class	Type	QC Passed	Location	Creation Date	Last Modified
<input type="checkbox"/>	SAM201216	PCSI_1073_Ly_R_nn_1-1_D_1	gDNA (aliquot)	GENOMIC	✓	6TH_LIBRARY_INBOX - B02 (647A_H_3_B-3_Lib-Inbox)	2019-05-24	2019-05-24 16:22:21
<input type="checkbox"/>	SAM201215	PCSI_1073_Lv_M_nn_1-1_D_1	gDNA (aliquot)	GENOMIC	✓	6TH_LIBRARY_INBOX - B01 (647A_H_3_B-3_Lib-Inbox)	2019-05-24	2019-05-24 16:22:20
<input type="checkbox"/>	SAM201214	PCSI_1069_Pa_P_nn_1-1_R_1	whole RNA (aliquot)	TRANSCRIPTOMIC	✓	RNA INBOX #2 - E02	2019-05-24	2019-05-24 16:21:19
<input type="checkbox"/>	SAM201213	PCSI_1064_Pa_P_nn_1-2_R_1	whole RNA (aliquot)	TRANSCRIPTOMIC	✓	RNA INBOX #2 - E01	2019-05-24	2019-05-24 16:21:18
<input type="checkbox"/>	SAM201209	SCRM_1040	Identity	TRANSCRIPTOMIC	?	Unknown		2019-05-24 16:02:02
<input type="checkbox"/>	SAM201210	SCRM_1040_Bn_X_nn_1-1	Tissue	TRANSCRIPTOMIC	?	Unknown		2019-05-24 16:02:02
<input type="checkbox"/>	SAM201211	SCRM_1040_Bn_X_nn_1-1_D_S1	cDNA (stock)	TRANSCRIPTOMIC	?	Unknown		2019-05-24 16:02:02
<input type="checkbox"/>	SAM201212	SCRM_1040_Bn_X_nn_1-1_D_1	cDNA (aliquot)	TRANSCRIPTOMIC	?	Unknown		2019-05-24 16:02:02

GOALS

1. Produce clean data

1. Make data entry not suck

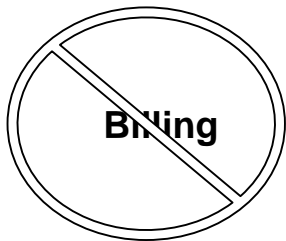
MISO Design

1. For who? Technicians
2. For what? Tracking laboratory activities
3. Play nicely with other software
4. Take tips from 40+ years of user-centred design

1. Designed for technicians

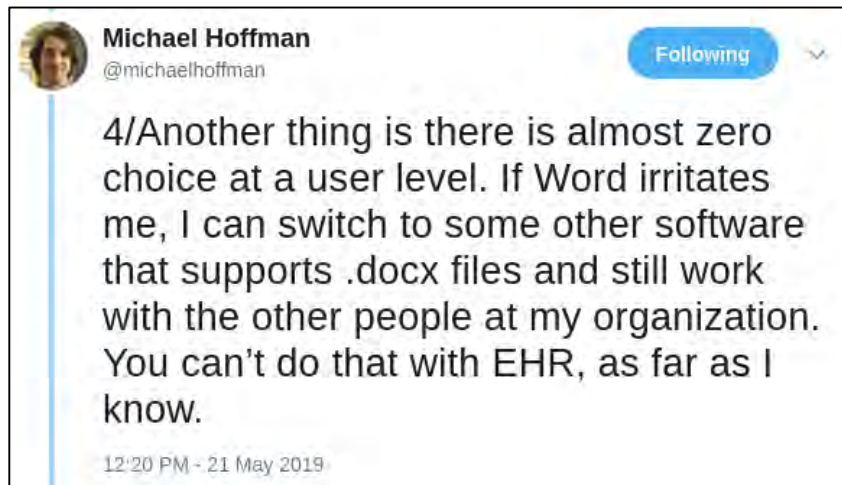


2. Specifically for laboratory tracking



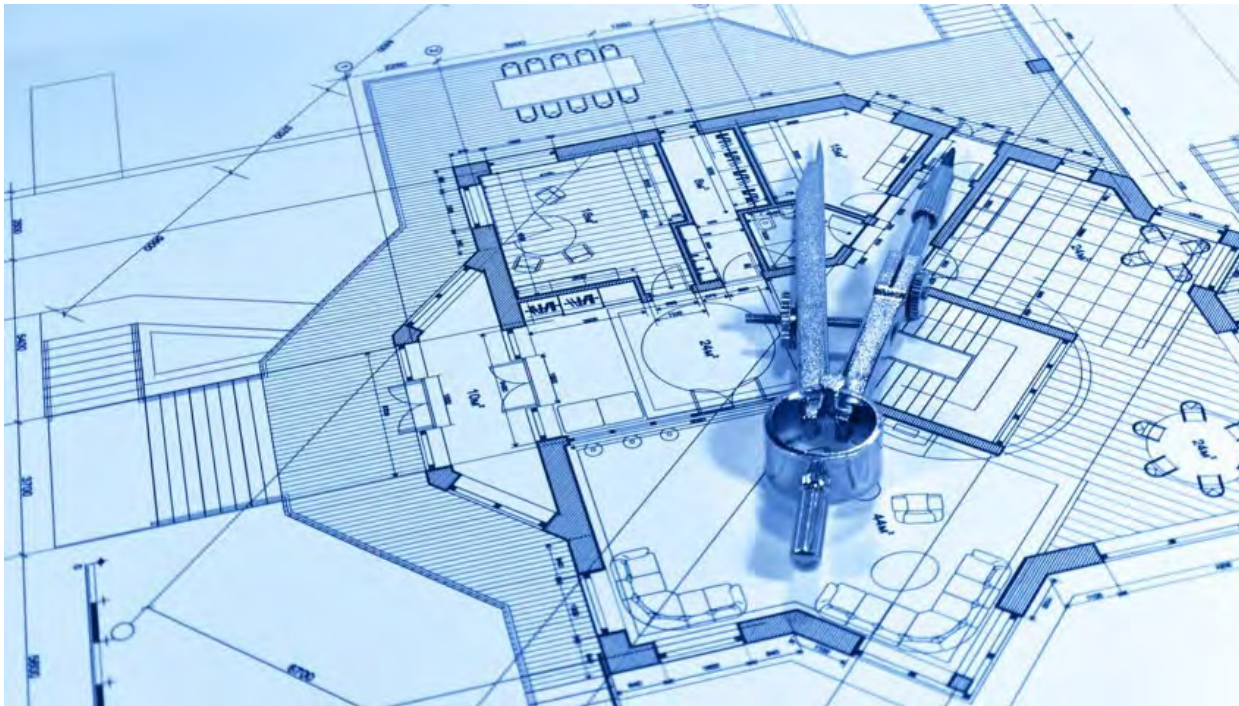
3. Play nicely with other software

- People are going to use Excel -- embrace it
- Let other people extend your software - provide application programming interfaces (APIs)

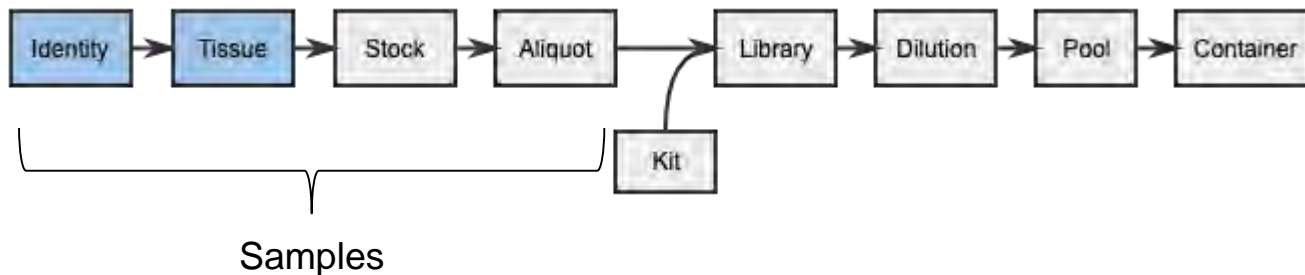


<https://twitter.com/michaelhoffman/status/1130870963275522049>

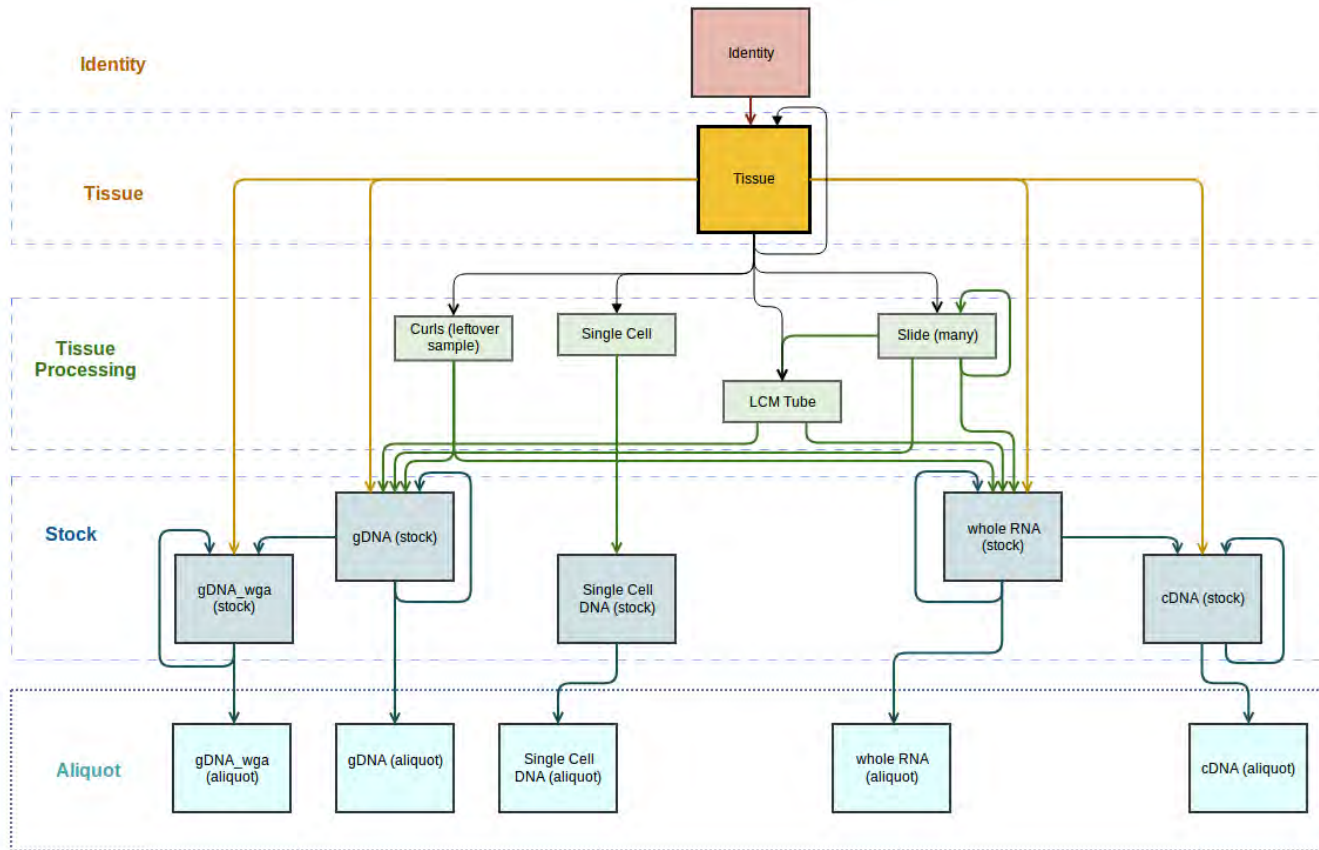
4. Tips from user interface design



Use the language of technicians

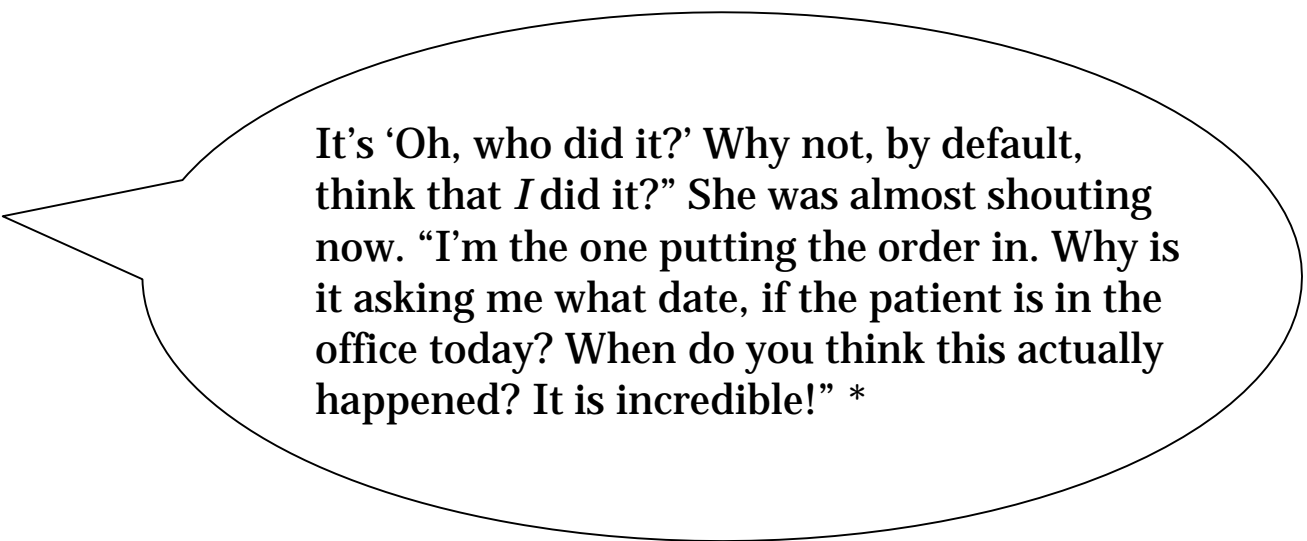


Change the software to fit the way they work



Reduce clicks

Doctor, while entering
electronic health
records



It's 'Oh, who did it?' Why not, by default, think that *I* did it?" She was almost shouting now. "I'm the one putting the order in. Why is it asking me what date, if the patient is in the office today? When do you think this actually happened? It is incredible!" *

1. Shortcuts to common tasks
2. Reasonable defaults (e.g. default to today's date)

* <https://www.newyorker.com/magazine/2018/11/12/why-doctors-hate-their-computers>

Use alerts *sparingly*



Don't overwhelm with information

Sort by Sample Location (rows)		Sort by Sample Location (columns)		Fill Boxes by Row	Fill Boxes by Column	Check QCs	Import	Export				
Library Name	Library Alias	Sample Alias	Sample Location	Matrix Barcode	Box Search	Box Alias	Position	Discarded	Distributed	Dist		
LIB27606	CSUR_0002_Bn_C_PE_370_CH	CSUR_0002_Bn_C_nn_1-1_D_1		0311450980		▼	▼	False ▼	Sent Out ▼	201'		
LIB27611	CSUR_0007_Hr_C_PE_373_CH	CSUR_0007_Hr_C_nn_1-1_D_1		0311450952		▼	▼	False ▼	Sent Out ▼	201'		
LIB30444	DCRT_015_Br_T_PE_232_EX	DCRT_015_Br_T_nn_1-1_D_1	6TH_LIBRARY_INBOX F05	0238515764		Seq_Inbox_NovaSeq (BOX799) ▼	H09 ▼	False ▼	No ▼			
LIB30445	DCRT_016_Br_P_PE_234_EX	DCRT_016_Br_P_nn_1-1_D_1	6TH_LIBRARY_INBOX F07	0311449220		Seq_Inbox_NovaSeq (BOX799) ▼	B10 ▼	False ▼	No ▼			
LIB30447	DCRT_016_Br_T_PE_270_EX	DCRT_016_Br_T_nn_1-1_D_1	6TH_LIBRARY_INBOX F08	0238515753		Seq_Inbox_NovaSeq (BOX799) ▼	G10 ▼	False ▼	No ▼			
LIB30448	DCRT_018_Br_P_PE_242_EX	DCRT_018_Br_P_nn_1-1_D_1	6TH_LIBRARY_INBOX F10	0125881713		Seq_Inbox_NovaSeq (BOX799) ▼	H10 ▼	False ▼	No ▼			
LIB30449	DCRT_018_Br_R_PE_255_EX	DCRT_018_Br_R_nn_1-1_D_1	6TH_LIBRARY_INBOX F12	0238515726		Seq_Inbox_NovaSeq (BOX799) ▼	A11 ▼	False ▼	No ▼			
LIB30450	DCRT_018_Br_T_PE_349_EX	DCRT_018_Br_T_nn_1-1_D_1	6TH_LIBRARY_INBOX F11	0311449203		Seq_Inbox_NovaSeq (BOX799) ▼	F11 ▼	False ▼	No ▼			
LIB30451	DCRT_020_Br_P_PE_231_EX	DCRT_020_Br_P_nn_1-1_D_1	6TH_LIBRARY_INBOX G10	0238515725		Seq_Inbox_NovaSeq (BOX799) ▼	G11 ▼	False ▼	No ▼			
LIB30452	DCRT_020_Br_R_PE_243_EX	DCRT_020_Br_R_nn_1-1_D_1	6TH_LIBRARY_INBOX G12	0311449013		Seq_Inbox_NovaSeq (BOX799) ▼	H11 ▼	False ▼	No ▼			
LIB30453	DCRT_020_Br_T_PE_245_EX	DCRT_020_Br_T_nn_1-1_D_1	6TH_LIBRARY_INBOX G11	0125882362		Seq_Inbox_NovaSeq (BOX799) ▼	B12 ▼	False ▼	No ▼			
LIB30454	GLCS_0001_Lv_R_PE_248_EX	GLCS_0001_Lv_R_nn_1-1_D_92		0238515758		Seq_Inbox_NovaSeq (BOX799) ▼	F12 ▼	False ▼	No ▼			
LIB30468	MNP_0004_Ad_M_PE_300_WG	MNP_0004_Ad_M_nn_1-1_D_1				▼	▼	False ▼	No ▼			
LIB30469	MNP_0005_Ad_M_PE_300_WG	MNP_0005_Ad_M_nn_1-2_D_1				▼	▼	False ▼	No ▼			
LIB30470	MNP_0006_Ad_M_PE_300_WG	MNP_0006_Ad_M_nn_1-3_D_1				▼	▼	False ▼	No ▼			
LIB30464	PCSI_1068_Lv_M_PE_388_WG	PCSI_1068_Lv_M_nn_1-1_D_1		0311449257		Seq_Inbox_NovaSeq (BOX799) ▼	A06 ▼	False ▼	No ▼			
LIB30465	PCSI_1068_Lv_M_PE_390_WG	PCSI_1068_Lv_M_nn_1-1_D_1		0311449287		Seq_Inbox_NovaSeq (BOX799) ▼	B07 ▼	False ▼	No ▼			

Plan for mistakes and let people know immediately

Pool Name	Pool Alias	Description	Matrix Barcode	Box Search	Box Alias	Position	Discarded	Creation Date	Concen
	MNP_POOL				▼	▼	False ▼	2019-05-26 ▼	

Missing Barcodes ✕

Pools should usually have barcodes. Are you sure you wish to save 1 pool without one?

Save **Cancel**

LISTEN TO YOUR USERS

Bug reports are gifts

Making data entry not suck

1. Plan specifically who and what the software is for (and stick to that)
2. Play nicely with other software
3. Use the language of your users
4. Fit the software to the protocol
5. Reduce clicks
6. Use alerts sparingly
7. Don't overwhelm with information
8. Plan for mistakes and give feedback immediately
9. Listen to your users

I didn't curse in my head even once!

- One of the greatest compliments we received on MISO

Think of the users



Think of the users



“A person-centred computer technology... can help foster a mature and humane society... If we can not enhance their working environment, what other good can we claim?”

Rob Kling, Towards a Person-Centered Computer Technology, 1973

Thank you

Acknowledgements:

- MISO team for their diligent work and excellent communication skills: Heather Armstrong, Dillan Cooke, Andre Masella, Alexis Varsava
- OICR Genomics, Diagnostic Development, and Translational Genomics Laboratory for their feedback, bug reports, and trust
- Lars Jorgensen and the rest of the Genome Sequence Informatics team





Funding provided by the
Government of Ontario.

