



The importance of the SCE in enabling our shift from proprietary programming to open-source data science

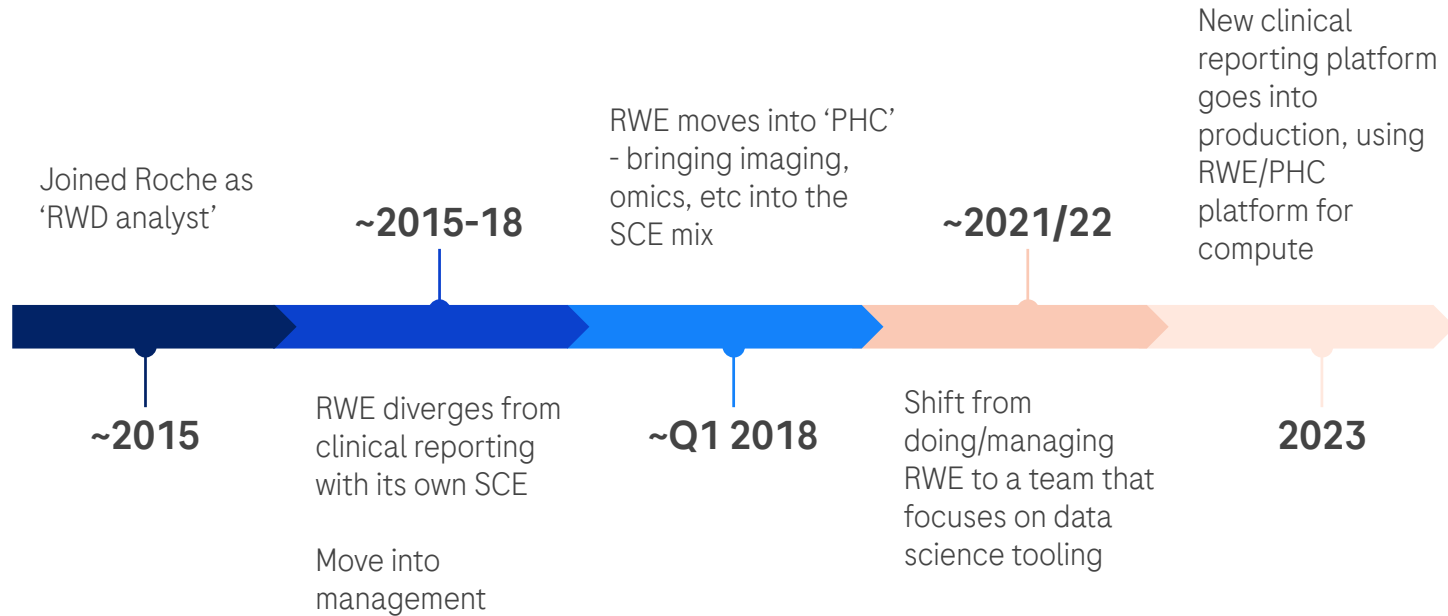
James Black, Data & Statistical Sciences, Roche

Talk contents

- My context
- RWE/PHC experiences
- What is an SCE?
- Bringing a modern SCE for data science into clinical reporting

My context

My journey with compute at Roche

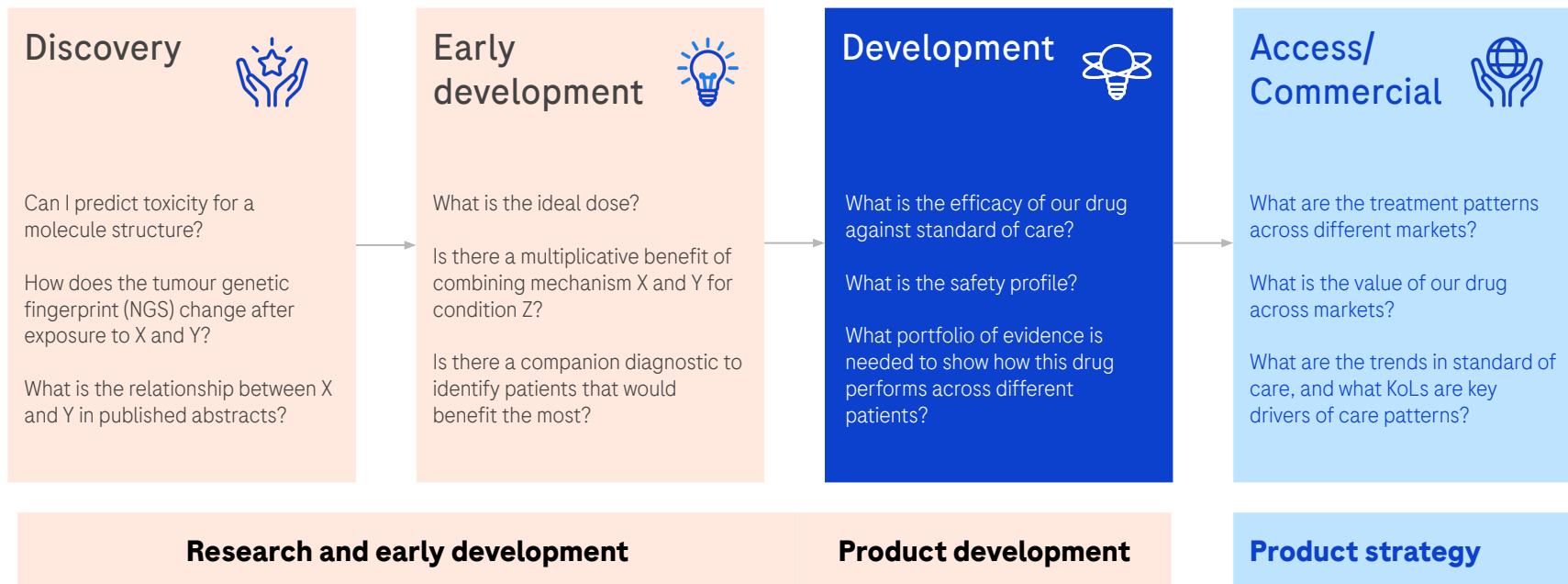


My current role's connection to SCEs



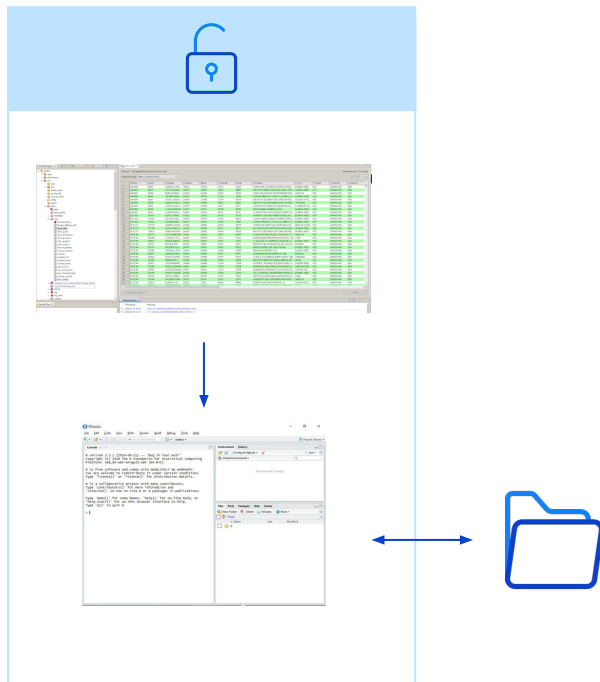
- Insights Engineering Product Family Lead
 - Sponsor our pan-study 'insights' codebase (e.g. Nest)
- Business lead for compute for late-stage
 - RWD/AA compute (Apollo) business owner
 - Clinical trial (Ocean) compute product owner
- Roche representative on the PHUSE SCE Council

Data science - and compute needs are diverse within a Pharma company



RWE/PHC experiences

RWE in 2015



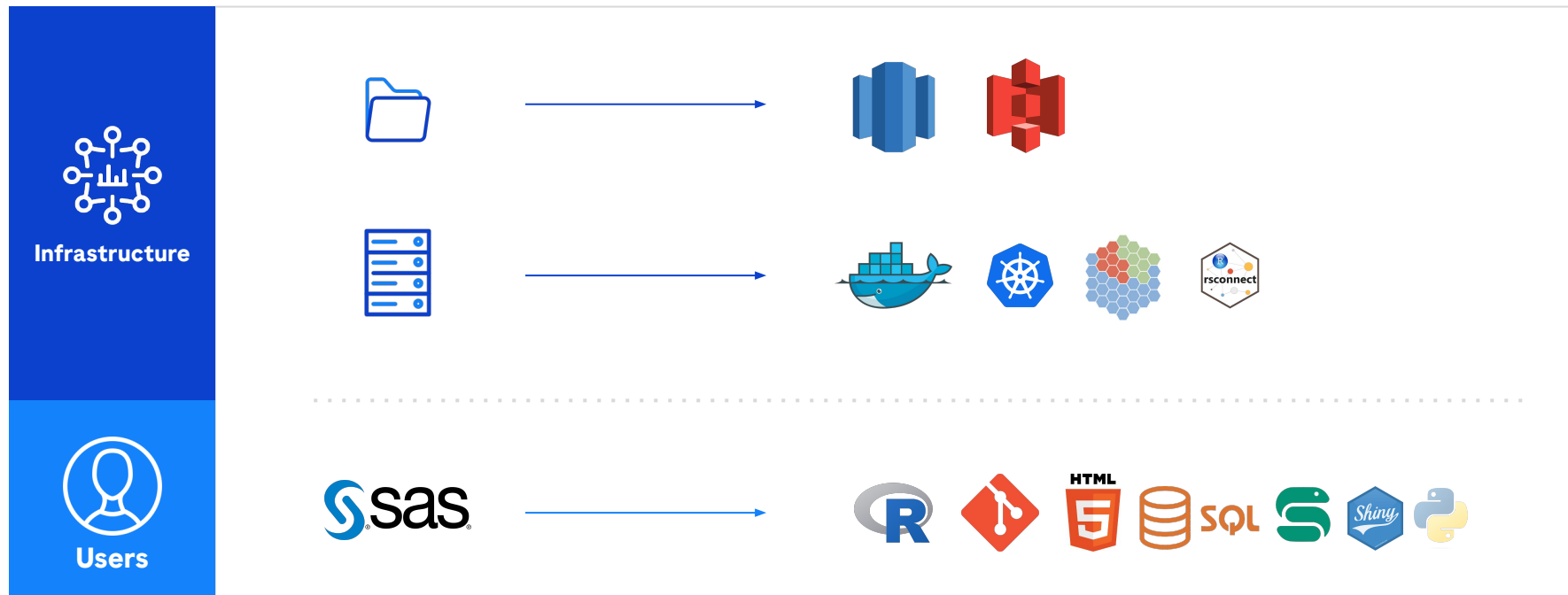
The environment

- Clinical trial data (and compute) in a SAS only ‘kitchen sink’ platform
- Data can be exported one-way into an R server
- Managed by IT with in-frequent releases
- RWD on a network drive (raw and derived)

The experience

- Reading data takes a long time, and data is duplicated as convenience copies
- Have to make tickets for system libraries
- Local laptops are a preferred place to work

RWE tooling before and after 2018



Processes and infrastructure catalysing data science adoption



- Git used for version control → Users familiar with github/gitlab and gitflows
- Snakemake for orchestration → Understanding of codebased workflow tooling
- Production runs off git repos → Familiar with CICD and 'don't trust interactively run code' mindset
- Redshift/S3 for data → Understanding of secrets management, APIs, and databases
- SAS to R → Exposed to huge open source libraries, external collaboration and thinking through dependencies
- Managed servers to containers → Users can take ownership of their environment at the system level

What is an SCE?

What is an SCE?

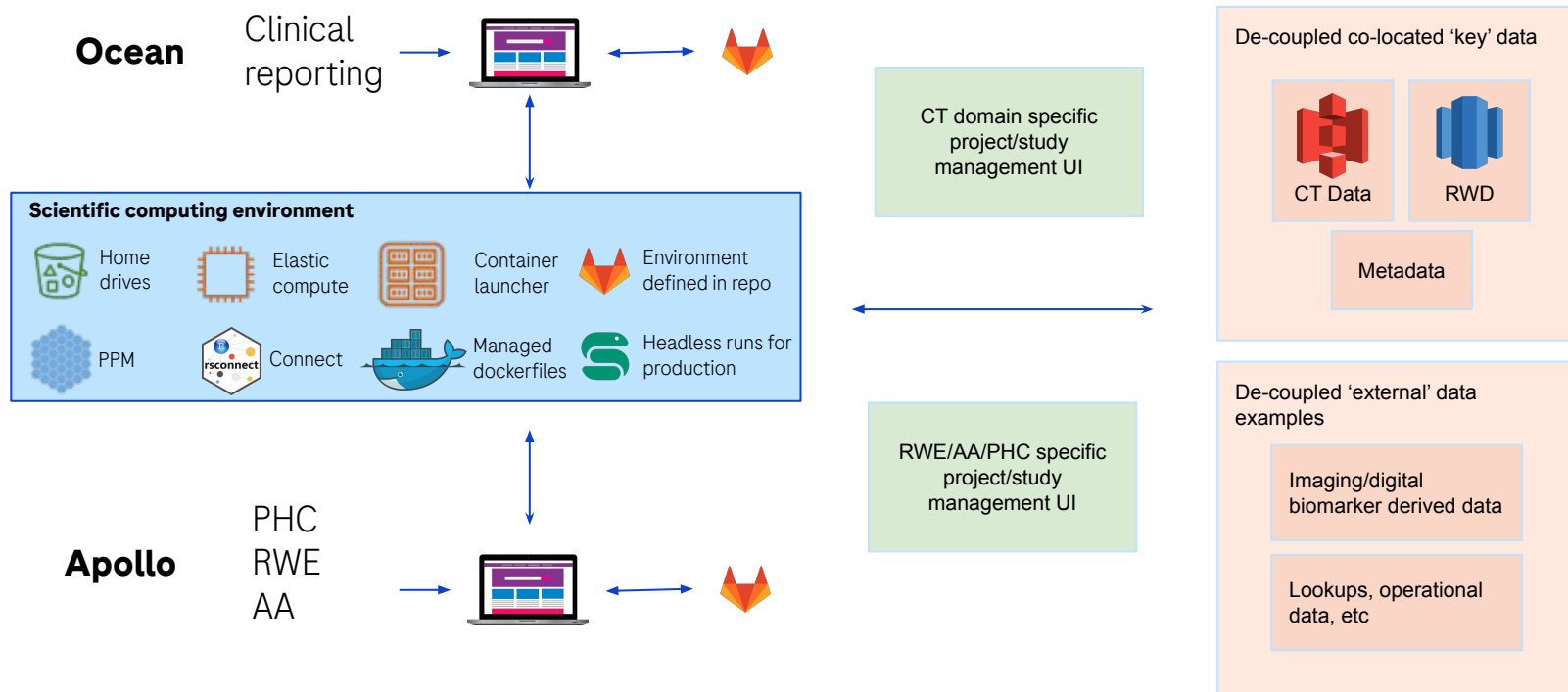


Statistical
Computing
Environment
White Paper

Authors: Mark Bynens (JnJ), Sheetal Patel (GSK), Olivier Leconte (JnJ), Delyth Jones (GSK), Sam Warden (GSK), Sascha Ahrweiler (Bayer), Oliver Richter (Boehringer Ingelheim), Jorine Putter (GSK), Joseph Rowley (Novartis), Marie-Claude Laramée (Novartis), Des Burke (GSK), Holger Dach (Bayer), Mario Lozina (Boehringer Ingelheim), Jon-Paul Mewes (Roche), Paul Fioole (JnJ), Eunice Ndungu (Merck), Mary Kuklinski (BMS), Timothy Kelly (BMS), Timothy Stuart Pearce (Pfizer) and Gary Chen (Pfizer).

- *Statistical or Scientific Computing Environment?*
- *“A modern SCE must allow for functional, meaningful and delightful experiences to the end users.”*
- *What is the scope?*
 - CDISC derived submission data and TLGs for GxP?
 - Exploratory work?
 - Trial design?
 - RWE?
 - Bioinformatics, imaging, digital biomarkers?
 - SaMD development/MLops?
- **Is data storage (source and derived) part on an SCE?**

What an SCE looks like at Roche

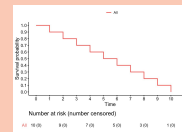


Validation shouldn't kill innovation in your SCE

Roche's approach to isolate validated insights from users normal workflows



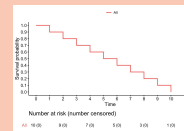
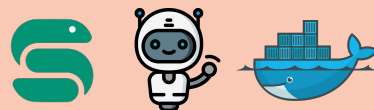
Work in IDE



	Age	Sex	Weight	Height	Weight	Height	Weight	Height	Weight	Height
Age	10	10	10	10	10	10	10	10	10	10
Sex	10	10	10	10	10	10	10	10	10	10
Weight	10	10	10	10	10	10	10	10	10	10
Height	10	10	10	10	10	10	10	10	10	10
Weight	10	10	10	10	10	10	10	10	10	10
Height	10	10	10	10	10	10	10	10	10	10
Weight	10	10	10	10	10	10	10	10	10	10
Height	10	10	10	10	10	10	10	10	10	10
Weight	10	10	10	10	10	10	10	10	10	10
Height	10	10	10	10	10	10	10	10	10	10



Batch runs



+



VALIDATED

Bringing a modern SCE for data science into clinical reporting

Scaling from RWE to clinical reporting



RWE: New-ish field with ~40 people touching data (in 2018, 500+ in 2023)



Clinical Trials: ~800 Statistical Programmers + additional Statisticians. 20+ years of tools and experience with prior processes

~~Scaling from RWE to clinical reporting~~

One simple trick to make a new SCE easy.....



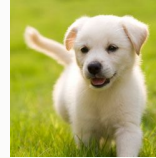
**Define how people work
from a blank slate**



**Support an evolution to a
new SCE and ways of
working**

Setting the dialogue on what a new SCE means to change management

“Existing processes will be supported as-is”



“A multilingual future”

Platform needs to support legacy
in an **efficient** way for the
foreseeable future

“Open source, with an R backbone”

Negotiated timelines to move to
new ways of working

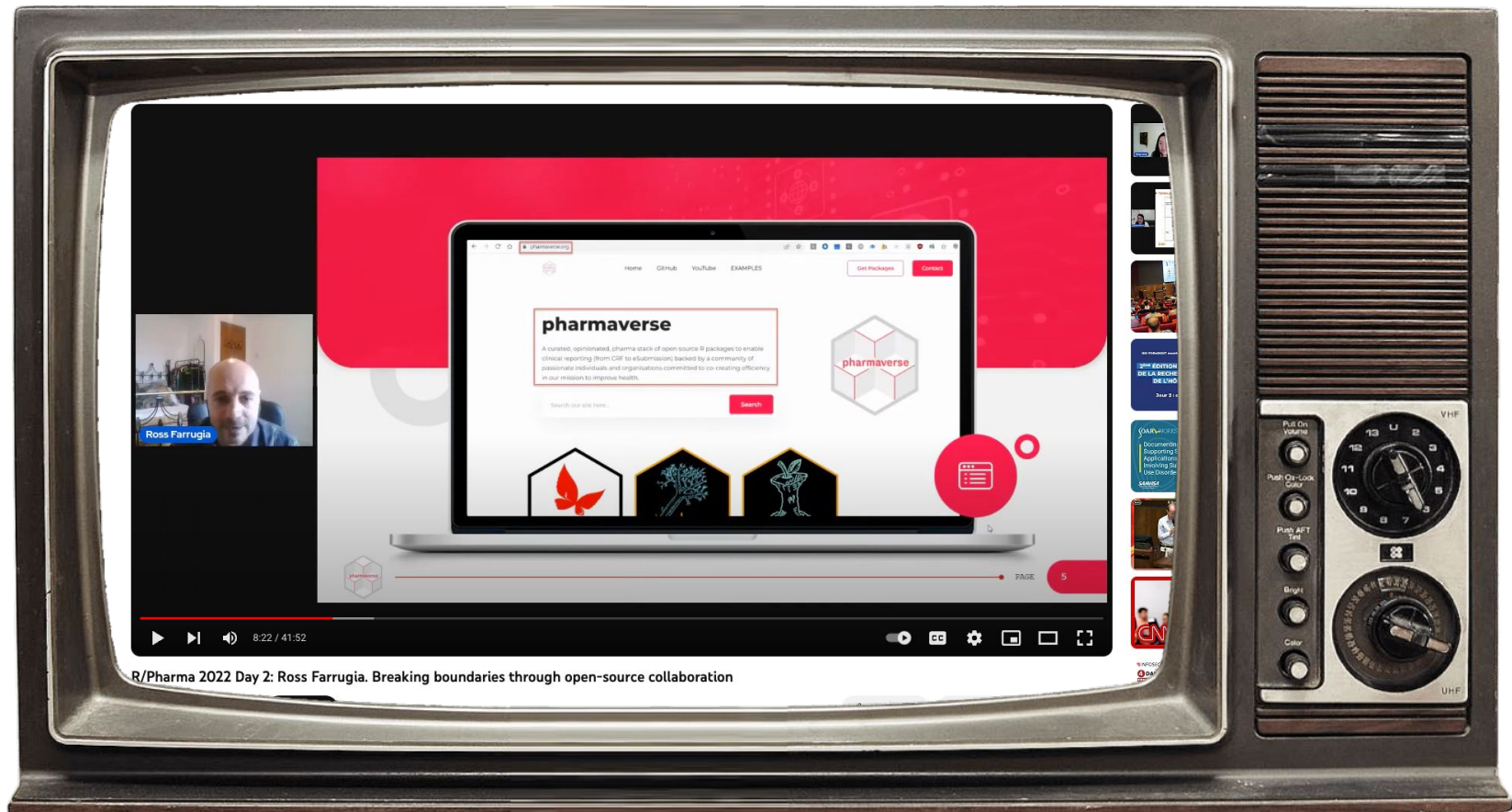
Why disrupt clinical reporting?



“We’ve spent decades with one language and vendor products that abstract away the technical aspects of data science into an end to end platform! Why throw that all away with a swamp of tools?”






“Git, codebased workflow management, sharing code via R packages - this all too technical and slows people down!”

“Statistical programming is a distinct, highly specialised and very process driven role - generic and flexible tools and platforms will only decrease efficiency”



R/Pharma 2022 Day 2: Ross Farrugia. Breaking boundaries through open-source collaboration

We need to refactor clinical reporting SCEs into more generalised data science SCEs

- Talent are graduating with experience in ,  and , so why can't statistical programmers use the same language as the statisticians?
- Clinical trial design and data modalities used continues to grow in complexity
- The  has shifted the codebase for clinical reporting to be a collaborative effort in 
 - Statistical programmers are now:
 - Spending more time in software development
 - Need to be able to handle new data modalities and emerging tools in other languages
 - Can be expected to have core data science competencies like git, environment management and keep up with data sciences evolution



Breaking clinical reporting out of its domain specific SCE design is required to support this shift to statistical reporting becoming a specialty within data science, rather than a separate silo'd world

Notes from R/Pharma round tables

“What are the bounds of an SCE?”

“We have many legacy workflows to support”

“Supporting legacy workflows (e.g. filesystem based data and old macros) cripples innovation”

“Maintaining provenance/replicability gets harder moving away from one platform”



Scene from the round tables in Chicago this year

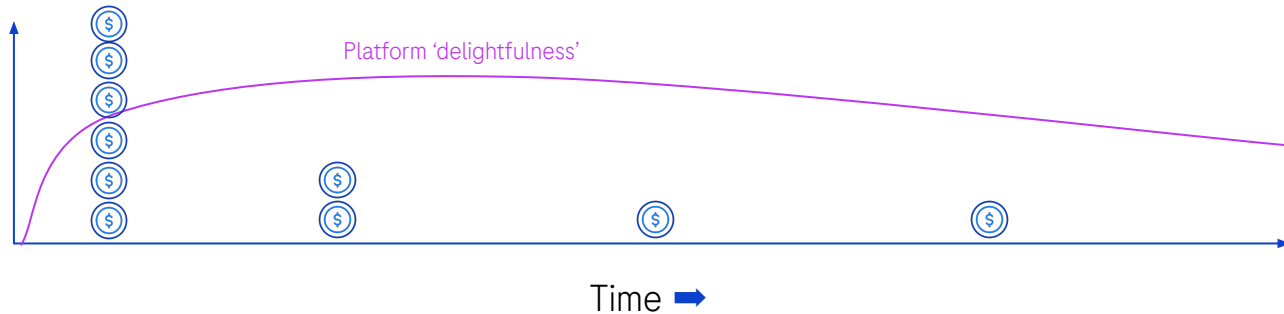
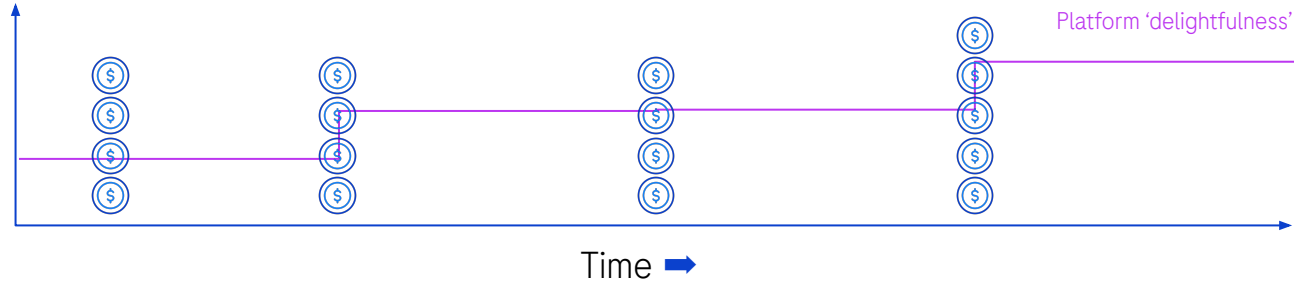
“How will a new SCE make the company money?”

“Can GxP and exploratory co-exist?”

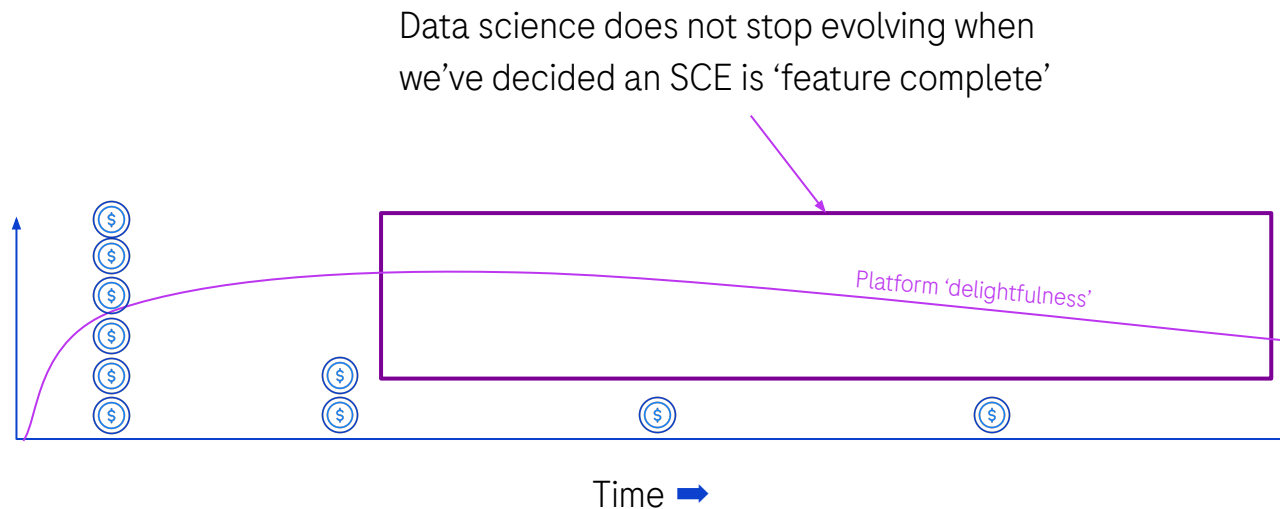
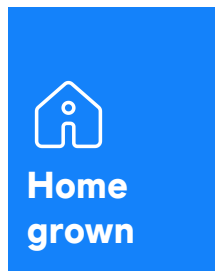
“Validation is painful and orientated towards vendor platforms”

“Sustained funding for homegrown post 1.0 is a challenge”

The curse of internal funding models



The curse of internal funding models



Re-capping some key questions to tackle with a modern clinical reporting SCE



- Is the scope clinical reporting by statistical programmers, late stage evidence, or even broader?
- Be clear on what you expect from users
 - Should statisticians use git/batchless runs for trial design?
 - Do we want statistical programmers on a platform optimised for specific tasks with lots of abstraction, or data scientists on an adaptive platform?
 - Can we isolate validation to a specific subset of the platform? (e.g. batch runs)
 - For what do we need to capture provenance/metadata? (e.g. batch runs)
 - Will/can legacy/existing studies be migrated?
 - If homegrown - can we find a sustainable budget for innovation?

The discussion continues later today....



Please also join us and the end of day 2 (today) for a panel discussion about what is a next-generation SCE



Rx Pharma

PANEL DISCUSSION

WHAT PUTS THE NEXT-GEN IN OUR NEXT GENERATION SCEs?

Mark Bynens (Johnson & Johnson)
Eileen Ching (GSK)
Pam Kalra (ZS)
Mary Kuklinski (Bristol Myers Squibb)
Kevin Kunzmann (Boehringer Ingelheim)
Eric Nantz (Eli Lilly)
Moderator: James Black (Roche)

**WED, OCT 25, 2023
1:40PM - 2:30PM EDT**

Register at <https://rinpharma.com/>

Doing now what patients need next