

# Employing Complex Datasets for More Effective Decision-Making in Drug Development

Fred Wilson  
Director, Clinical Imaging  
Experimental Medicine Imaging

[fred.2.wilson@gsk.com](mailto:fred.2.wilson@gsk.com)

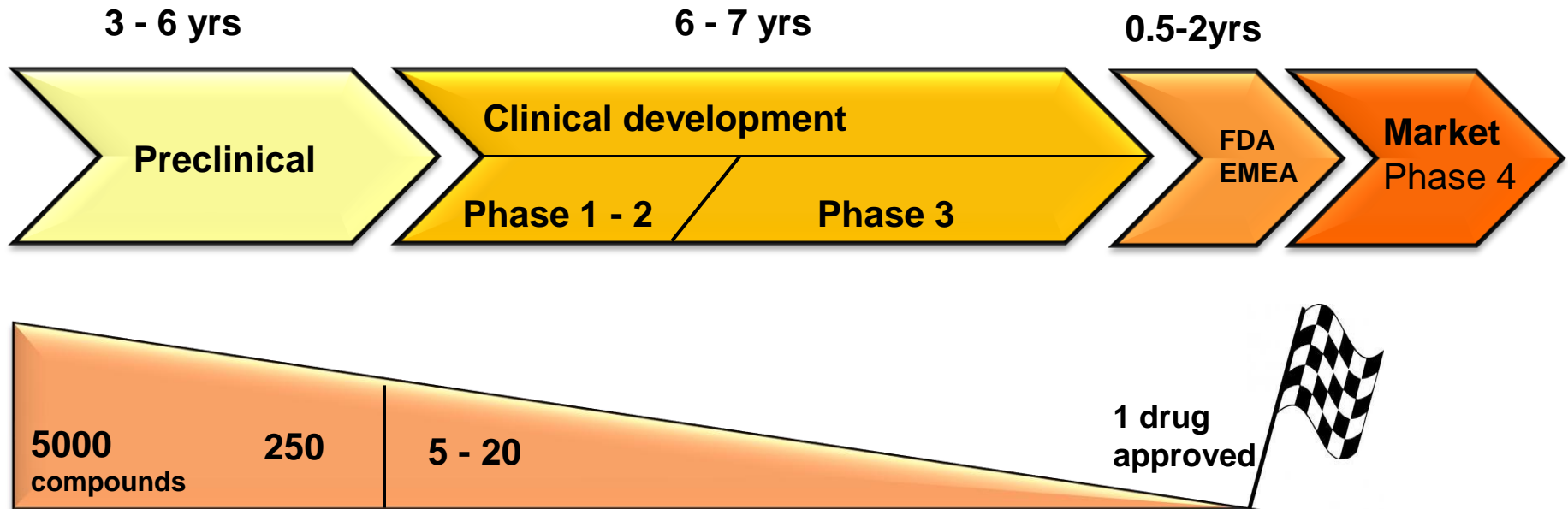
Chris Page  
Manager, Support Analyst  
Digital Delivery and Imaging

[chris.s.page@gsk.com](mailto:chris.s.page@gsk.com)

- Both presenters:
    - Current employees of GlaxoSmithKline and hold stock
  
  - Fred Wilson:
    - Previously a consultant to ECNP R&S, GlaxoSmithKline, IPPEC, King's College London, Lundbeck A/S, Mentis Cura ehf and Pfizer Inc.
    - Received travel expenses as a guest speaker on EEG from Orion Pharma Ltd
    - Previously an employee of Pfizer and held stock options
-

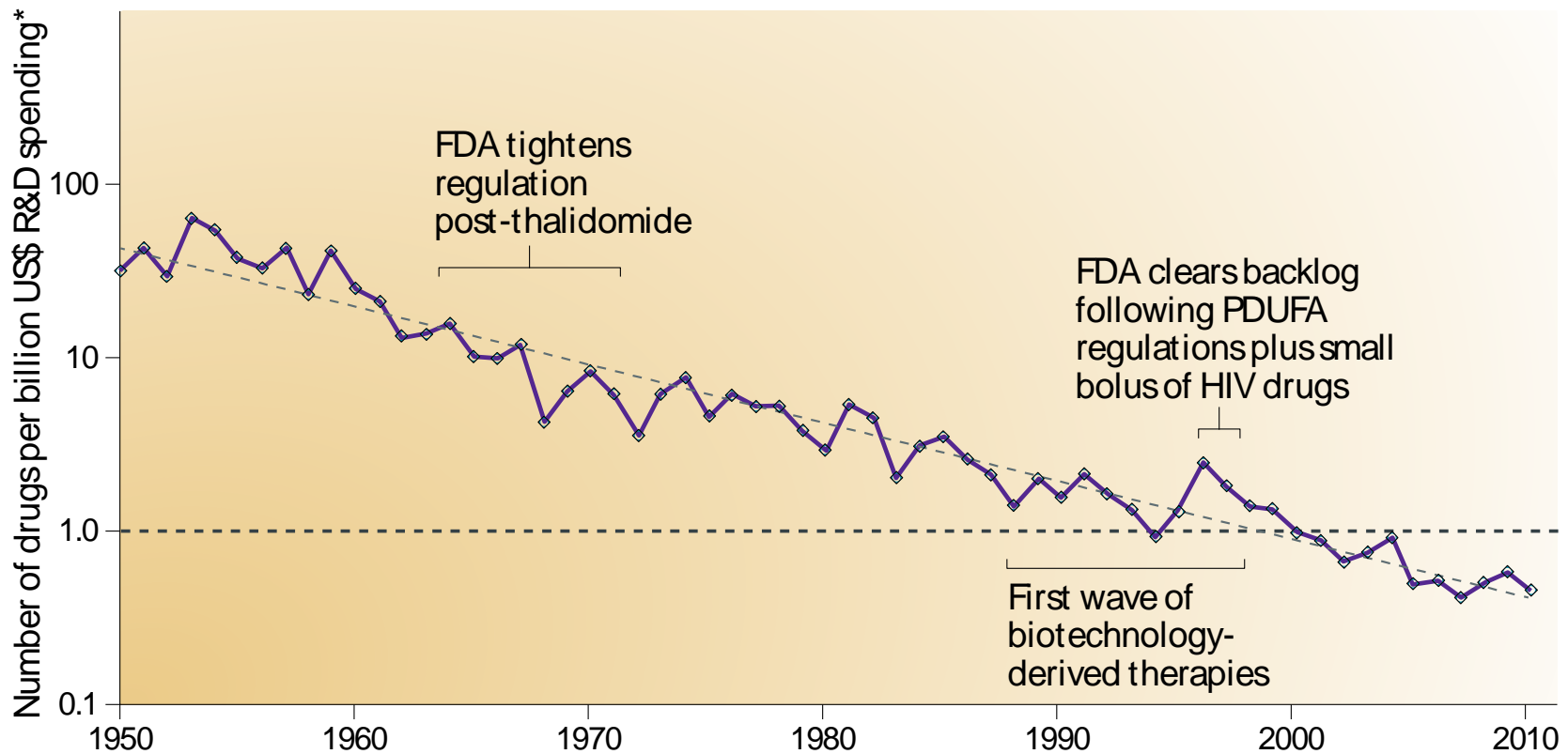
- 
- Motivation:
    - Attrition in the drug development pipeline
    - What do we mean by complex data and decision-making?
  - Improving decision-making in early drug development:
    - The role of biomarkers – what do we need to measure?
    - Example: electroencephalography (EEG) as a pharmacodynamic biomarker
  - Quality control and data linkage in multi-site clinical studies:
    - Improving on existing visual and other basic measures
    - Extracting additional information from existing datasets
  - Conclusions
-

# Drug development – a lengthy and “risky” process

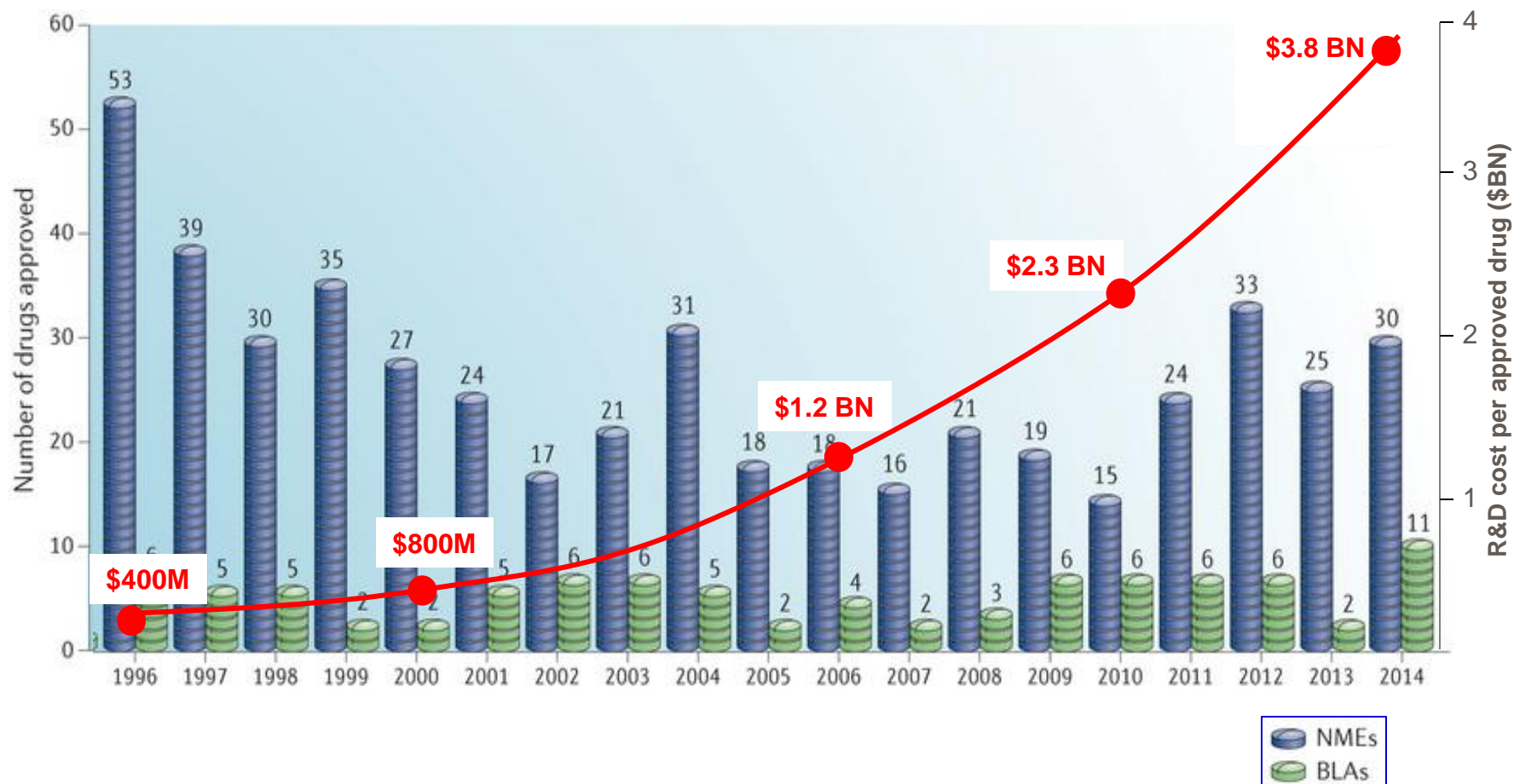


- The situation is economically unsustainable

## a Overall trend in R&D efficiency (inflation-adjusted)



# R&D productivity - increasing costs per approved drug



# What do we mean by complex data and decision-making?

---



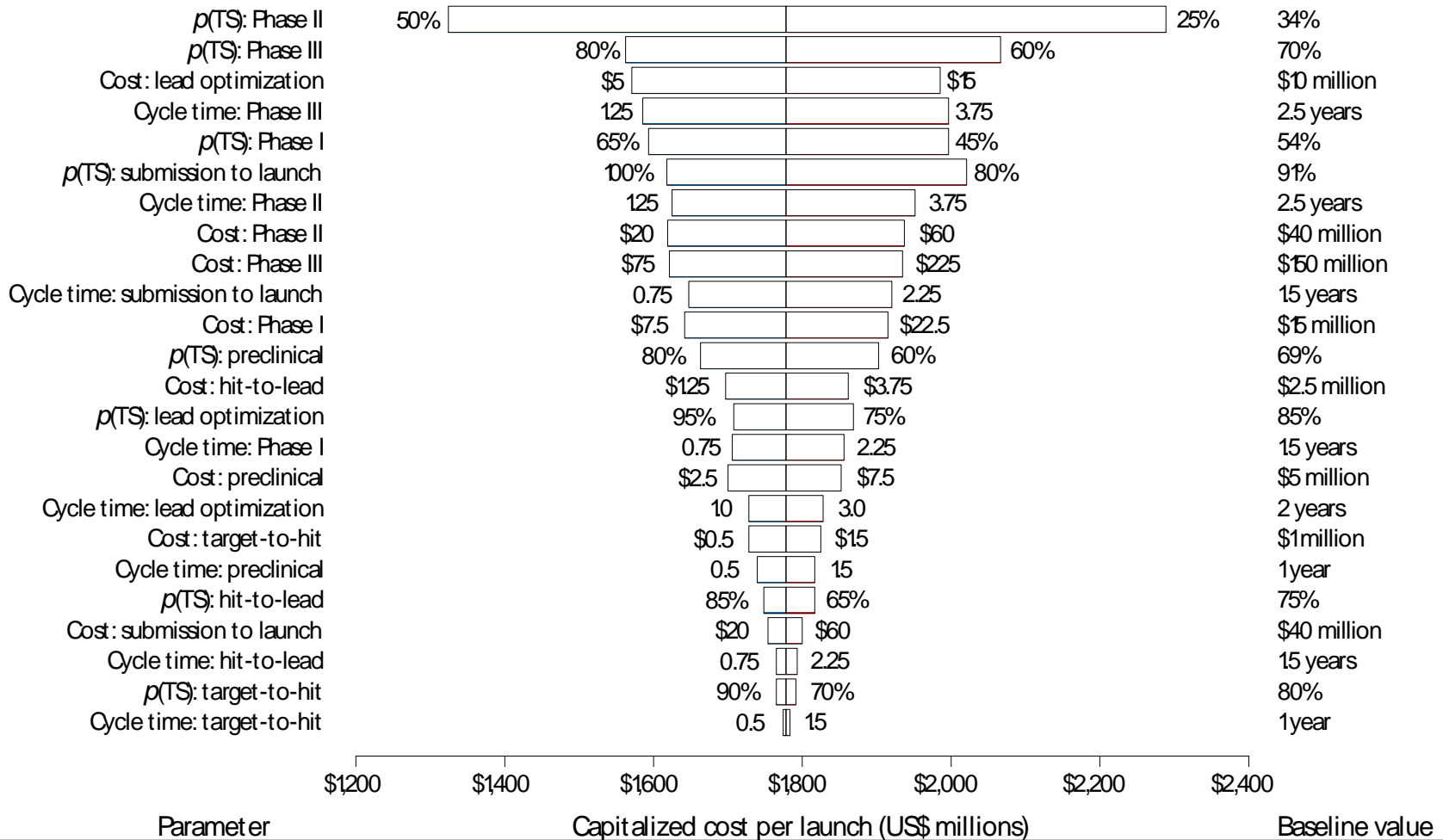
- Most biological and clinical datasets are ‘complex’:
  - Large numbers of data points
  - Multiple sources of noise (random, biological, systematic)
  - May not include large numbers of samples (so not true ‘big data’)
- ‘Decision-making’ requires data reduction to answer a specific question:
  - Typically requires a binary choice and/or reduction to a single variable, for example:
    - Is the drug binding to the target?
    - Is the drug having a biological effect? How big an effect?
    - Will this patient respond to the drug? By how much?



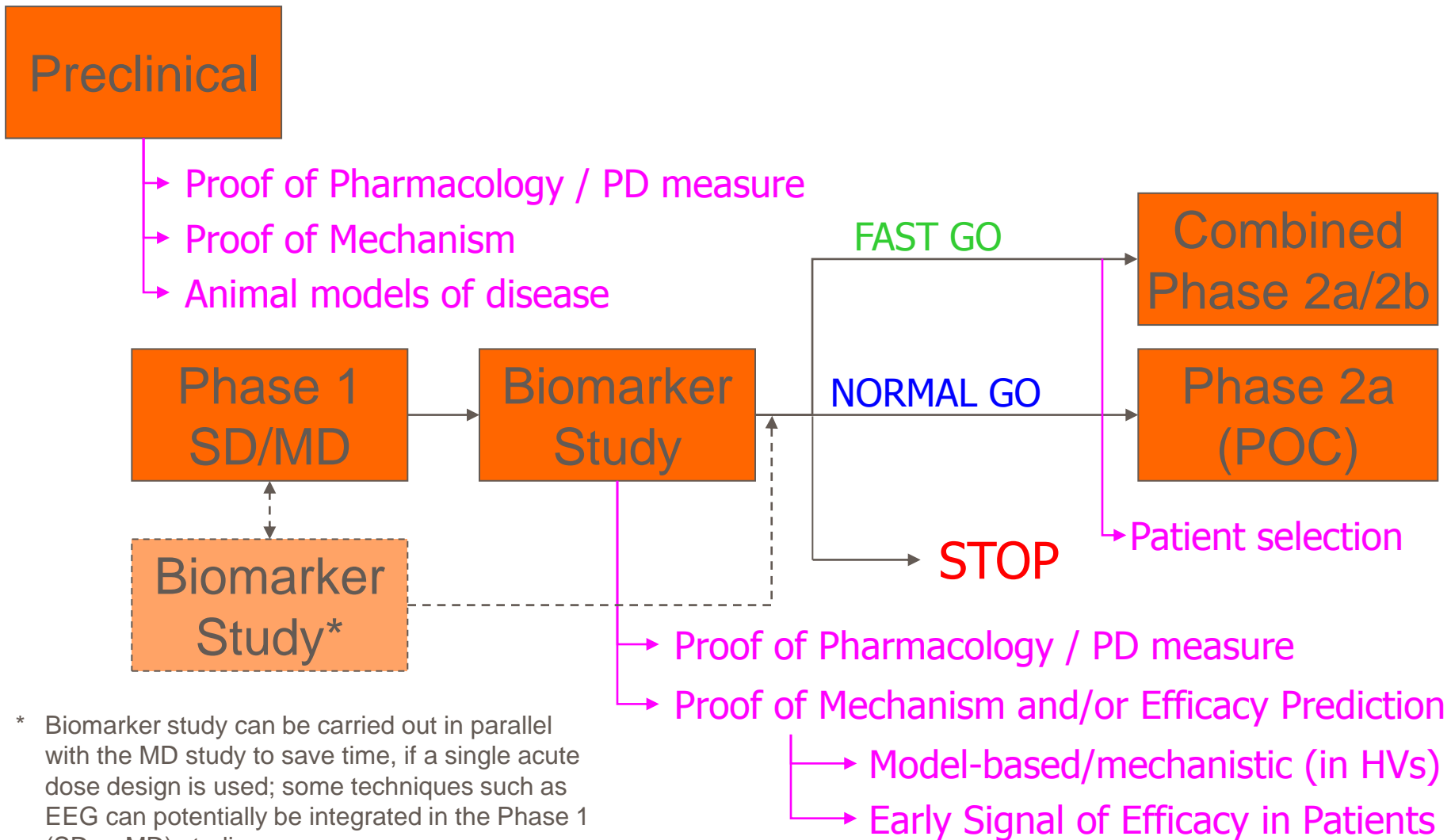
# Improving decision-making in early drug development



# Parametric Sensitivity Analysis



# De-risking Phase 2/3 using Biomarkers



\* Biomarker study can be carried out in parallel with the MD study to save time, if a single acute dose design is used; some techniques such as EEG can potentially be integrated in the Phase 1 (SD or MD) studies;

– Recent review of 44 Phase 2 drug development projects at Pfizer

– Examined based on 3 principles:

**PILLAR 1:** Exposure at the target site of action

**PILLAR 2:** Binding to the pharmacological target

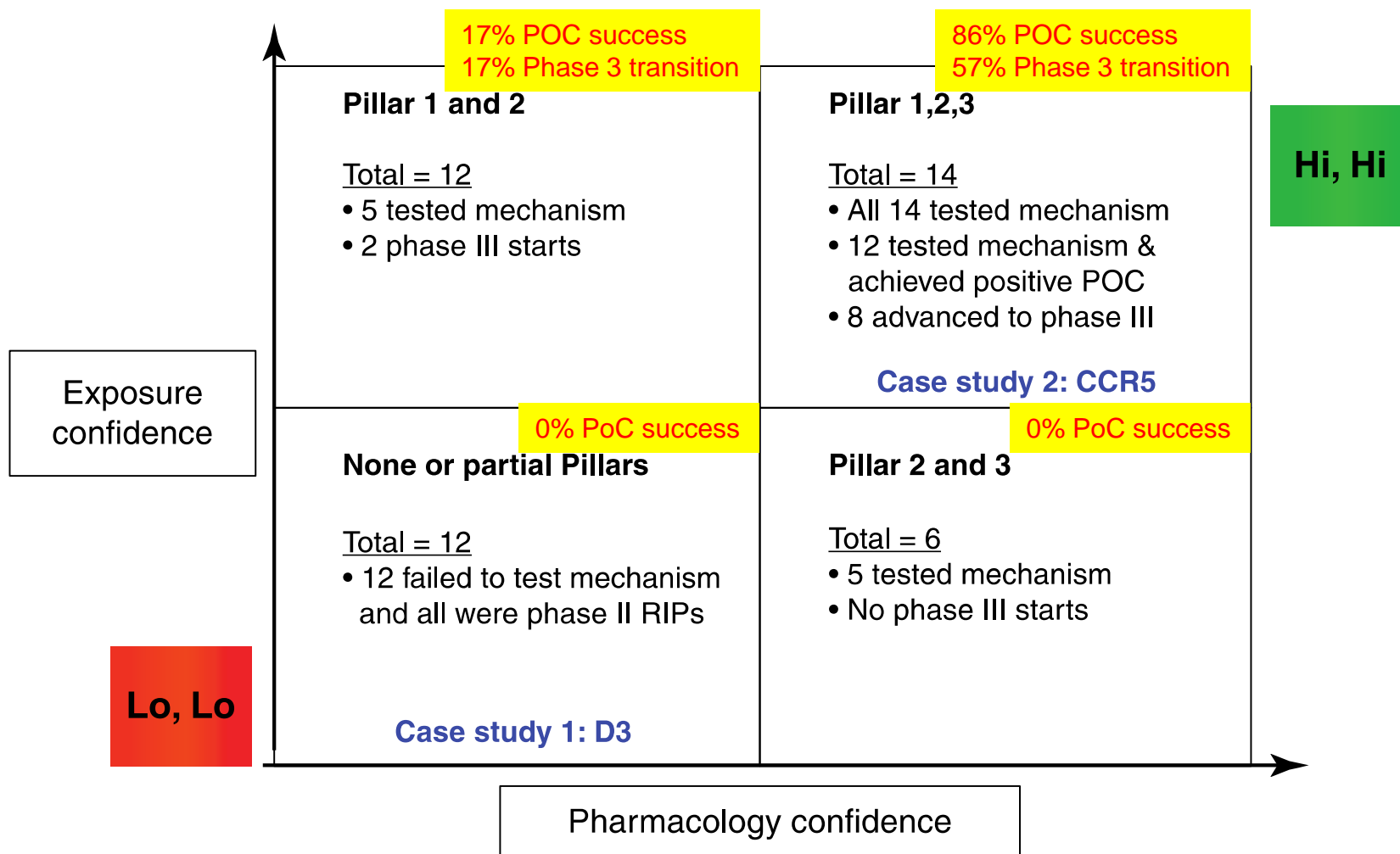
**PILLAR 3:** Expression of pharmacology

– Summarised onto two axes:

**EXPOSURE CONFIDENCE:** Based on Pillars 1 and 2

**PHARMACOLOGY CONFIDENCE:** Based on Pillars 2 and 3

# Fundamental PK-PD Principles



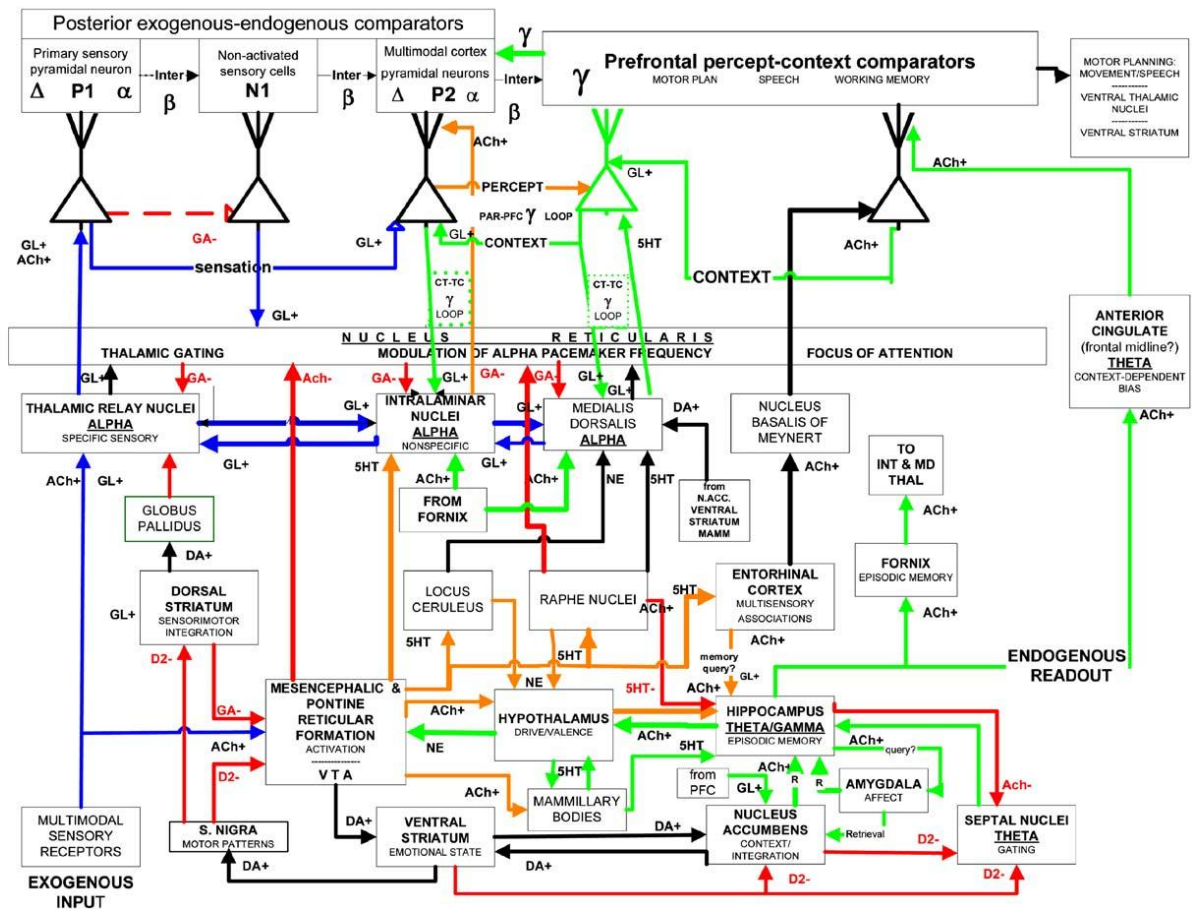
# EEG – a window onto brain function



Involves large neuronal populations that include all major neurotransmitter systems

## HOMEOSTATIC EEG REGULATORY SYSTEM

BLUE= EXOGENOUS SPECIFIC INPUT GOLD = NONSPECIFIC PROCESSING GREEN = ENDOGENOUS READOUT RED= INHIBITORY INFLUENCES

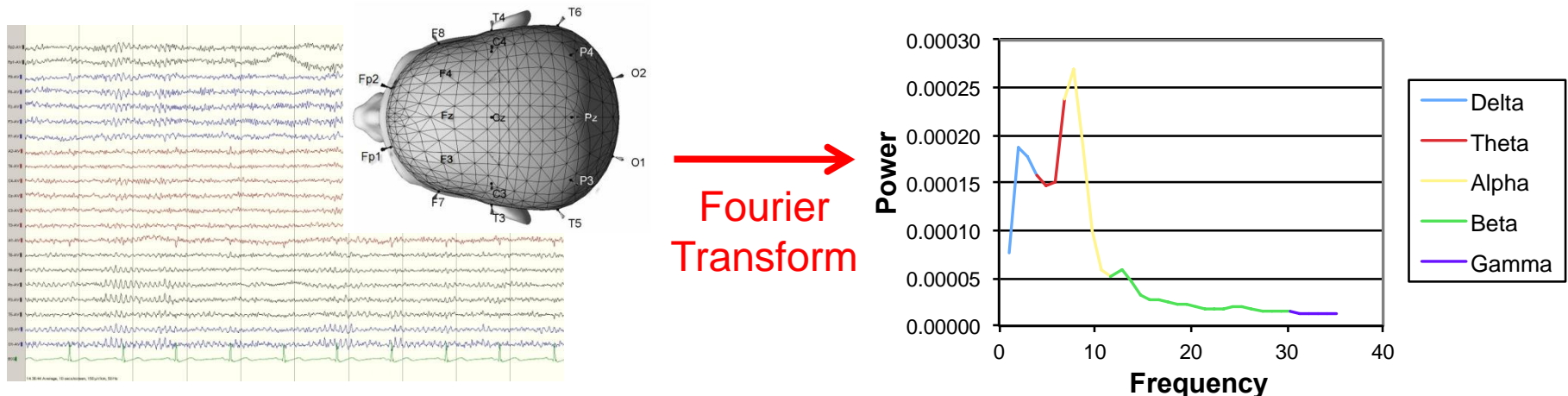


- 
- Lots of historical issues with unclear results from pEEG
  - Propose a new framework for when to use pEEG as a PD biomarker:
    - Two simple criteria:
      - Preclinical experiments produce a robust result
      - We expect this to translate (based on best current knowledge)
    - Clinical study should be designed to test for the expected effect, with other pEEG measures as secondary endpoints

# Classical Quantitative EEG Analysis



- Generate frequency spectrum of signal (e.g. using Short-Term Fast Fourier Transform)
- Split frequencies into bands (Delta, Theta, Alpha, Beta, Gamma)
- Evaluate required endpoints:
  - Total and relative spectral power in each band
  - Power ratios
  - Coherence between different regions in each frequency band
  - Other parameters e.g. peak alpha frequency

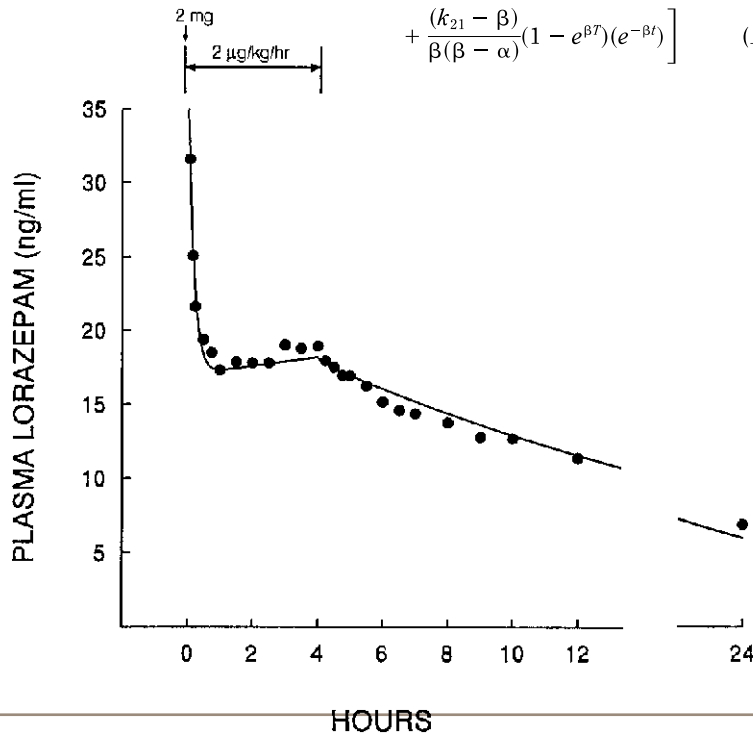


# Famous Example - Benzodiazepines



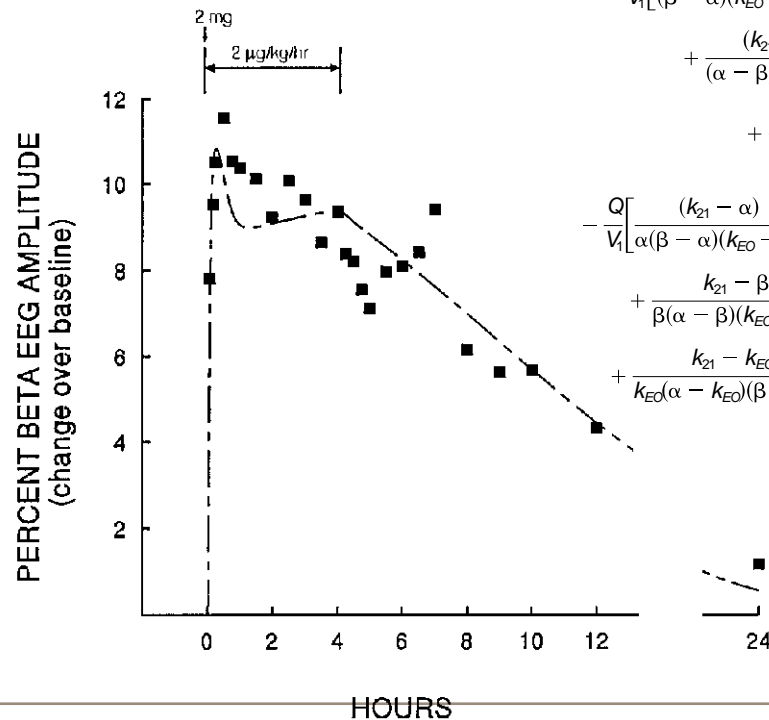
- Complex PK-PD modelling with EEG works well e.g.

$$C = \frac{D}{V_1} \left[ \frac{k_{21} - \alpha}{\beta - \alpha} \cdot e^{-\alpha t} + \frac{k_{21} - \beta}{\alpha - \beta} \cdot e^{-\beta t} \right] + \frac{Q}{V_1} \left[ \frac{(k_{21} - \alpha)}{\alpha(\alpha - \beta)} (1 - e^{\alpha T})(e^{-\alpha t}) + \frac{(k_{21} - \beta)}{\beta(\beta - \alpha)} (1 - e^{\beta T})(e^{-\beta t}) \right] \quad (1)$$



$$E = \frac{E_{max} \cdot (C_E \cdot k_{EO})^A}{EC_{50}^A + (C_E \cdot k_{EO})^A} \quad (2)$$

$$C_E = \frac{D}{V_1} \left[ \frac{(k_{21} - \alpha)}{(\beta - \alpha)(k_{EO} - \alpha)} \cdot e^{-\alpha t} + \frac{(k_{21} - \beta)}{(\alpha - \beta)(k_{EO} - \beta)} \cdot e^{-\beta t} + \frac{(k_{21} - k_{EO})}{(\alpha - k_{EO})(\beta - k_{EO})} \cdot e^{-k_{EO} t} \right] - \frac{Q}{V_1} \left[ \frac{(k_{21} - \alpha)}{\alpha(\beta - \alpha)(k_{EO} - \alpha)} \cdot (1 - e^{\alpha T})(e^{-\alpha t}) + \frac{k_{21} - \beta}{\beta(\alpha - \beta)(k_{EO} - \beta)} \cdot (1 - e^{\beta T})(e^{-\beta t}) + \frac{k_{21} - k_{EO}}{k_{EO}(\alpha - k_{EO})(\beta - k_{EO})} \cdot (1 - e^{k_{EO} T})(e^{-k_{EO} t}) \right] \quad (3)$$





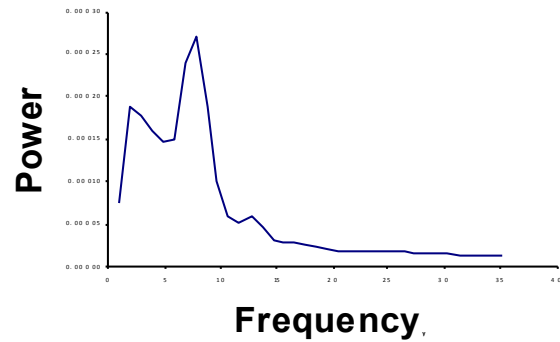
- 
- Numerous potential endpoints (100s or 1000s):
    - 19 or more electrode positions
    - 5 frequency bands (more if subdivided)
    - Absolute and relative power values
    - Power ratios
    - Coherence measures (by scalp region and band)
  - Individual endpoints lack specificity
  - Readout often dependent on *post hoc* interpretation
  - Impossible to define criteria *a priori* to enable clear decisions

# Generalised Semi-linear Canonical Correlation Analysis (GSLCCA)

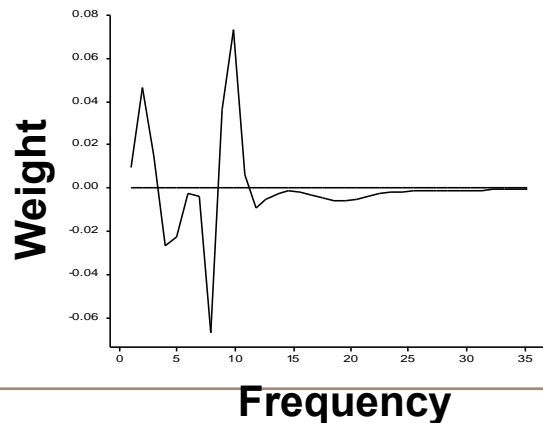


- Method developed to enhance utility of EEG as a PD biomarker by using data from the:
  - Whole spectrum (without dividing into bands)
  - Entire recording duration
  - All electrodes
- To provide:
  - Interpretable mechanistic information
  - A PD measure
- Assuming:
  - A PD profile of a known form (i.e. a given equation with unknown parameters)

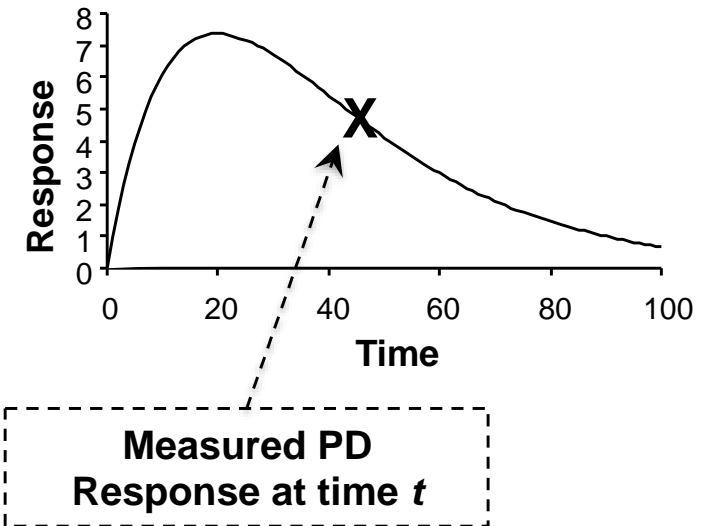
(a) Power spectrum at time  $t$



(b) Signature obtained using GSLCCA



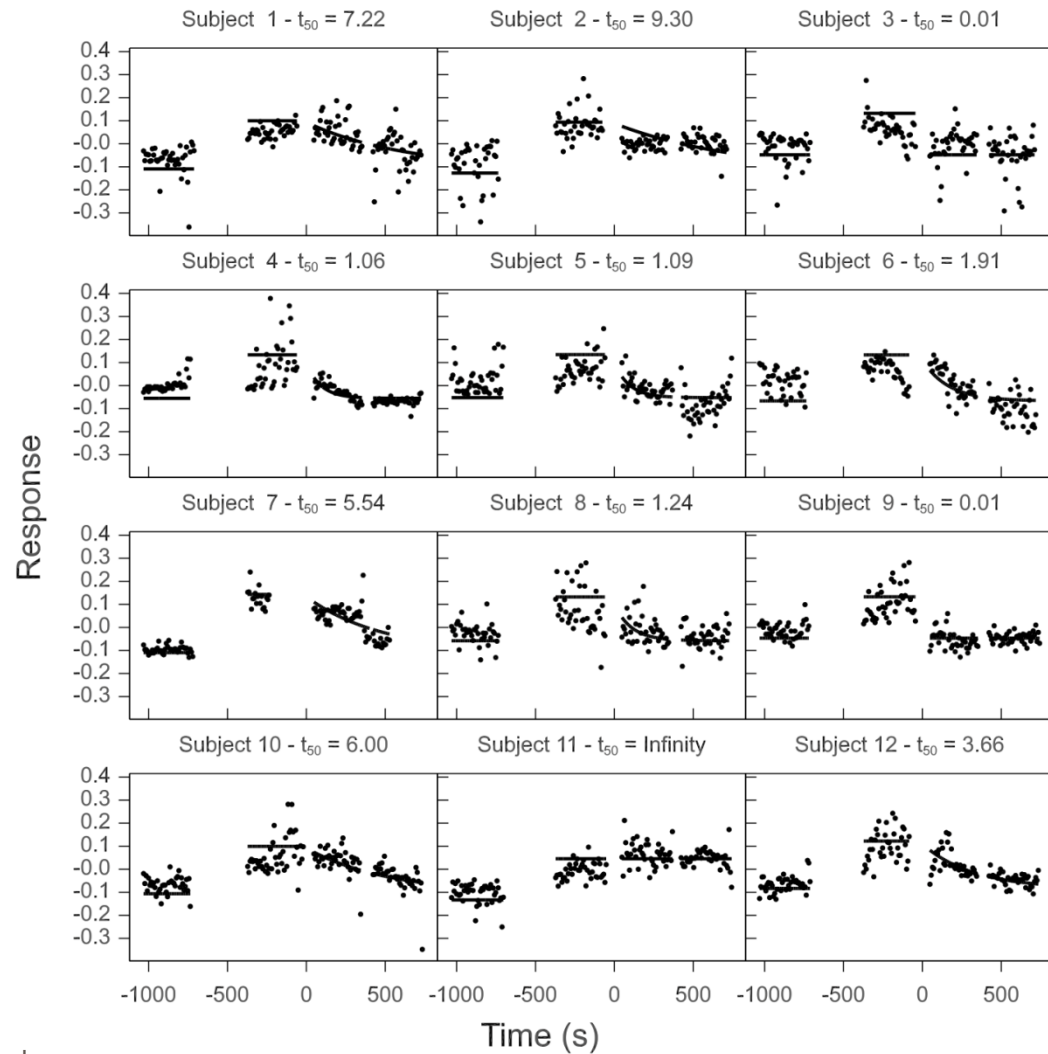
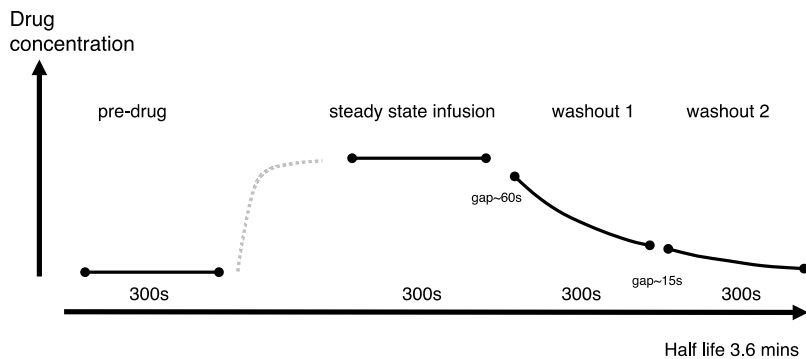
(c) Model PD response profile



# GSLCCA – Example Results



## Clinical study with remifentanyl



Mean  $t_{50} = 3.04 \pm 0.88$  minutes



# Quality control and data linkage in multi-site clinical studies

# Linking Imaging to Other Clinical Endpoints



## *Strategy for “Big Data” and Stratified Medicine*

---

**Goal of stratified medicine is to allow a clinician to determine the optimal therapy or combination of therapies for an individual at the earliest possible stage**

- How can this be determined based on initial presentation of disease?
  - Integrated analysis of genomic and other data
- Imaging is primary endpoint in many clinical studies
- Incorporating imaging data to analysis is challenging
  - Raw data are essentially large volumes of pixel intensities
  - Requires semantically-rich descriptors to correlate with other data sources
  - Essentially a problem of knowledge extraction from image volumes
- Not a classical Big Data problem
  - Relatively small number of samples (subject-visits)
  - Each sample is very well-characterised

# Registration-Path Imaging Studies



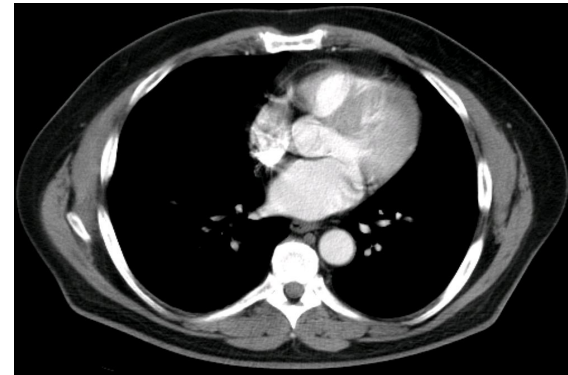
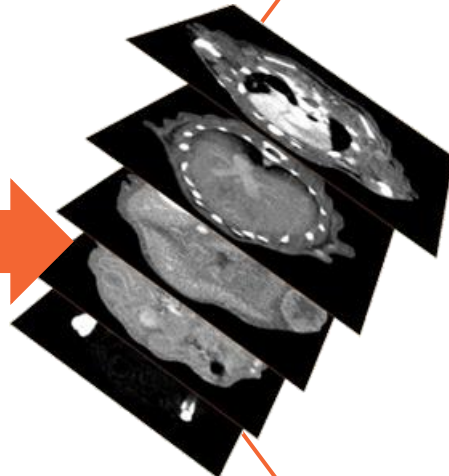
*Multisite and standardised*



- Safety and efficacy
- Established endpoints
  
- Large(ish) sample populations
- Data acquired globally in clinical radiology departments
- Local and centralised independent radiological review
  
- Regulated
- Conservative

# Clinical Imaging Data

*Digital Imaging and Communications in Medicine (DICOM)*



Image

Patient's Name	001234
Patient ID	001234
Patient's Sex	Male
Study Date	05-Dec-2007
Patient's Birth Date	30-Jun-1960
Modality	CT
Referring Physician	

Metadata



# Sensitive Personally Identifiable Information



## Pixel deidentification

21-Sep-2011 09:51

Station:  
Untersuchender Arzt: **[REDACTED]**  
Assistenz: **[REDACTED]**

Gesamt mAs 5128 Gesamt DLP 962

	Scan	kV	mAs / ref.	CTDIvol	DLP	TI	cSL
Patientenposition H-SP							
Topogramm	1	120				5.3	0.6
Ob.Bauch nativ	2	120	123 / 180	8.33	244	0.5	1.2
Premonitoring	3	120	20	4.04	4	0.5	1.2
I.V. Bolus							
Monitoring	4	120	20	28.31	27	0.5	1.2
arteriell	11	120	121 / 180	8.21	226	0.5	1.2
portalvenoes	12	120	141 / 180	9.55	461	0.5	1.2



# Response Evaluation Criteria in Solid Tumours (RECIST)

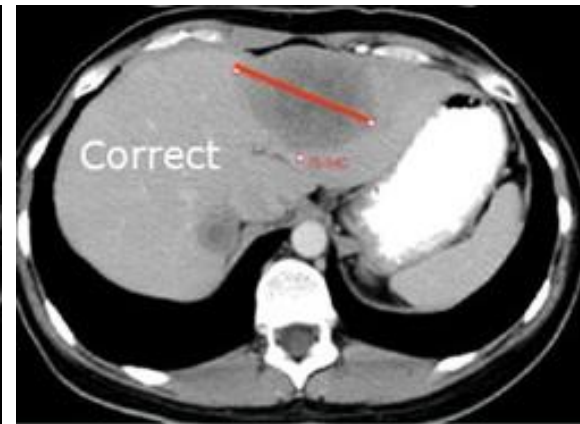
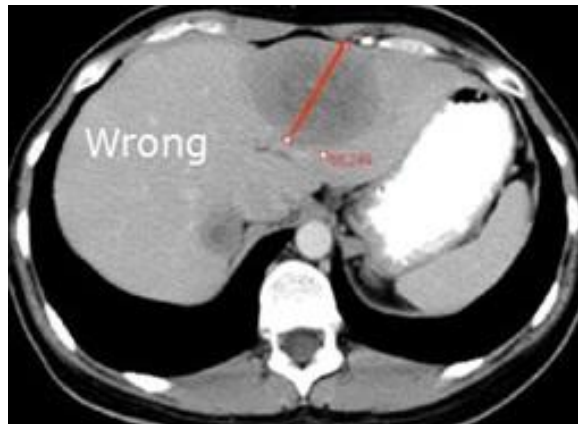


*Standard objective measures of response to therapy*

Baseline



Follow-on



<http://www.recist.com/recist-in-practice/>

# QC and Analysis Pipeline

*Opportunities for automation*



- Algorithms should be general
  - Validation overhead obviates study-specific software
  - Broad applicability across TAs
- Outputs should include confidence estimate
  - Need to be able to identify false-positives
- Challenges
  - Statistical bias: value of comparing data between studies?
  - Variations in acquisition (multisite)

# Classification and Automated QC



## *Randomised Decision Forests*



Courtesy Ben Glocker

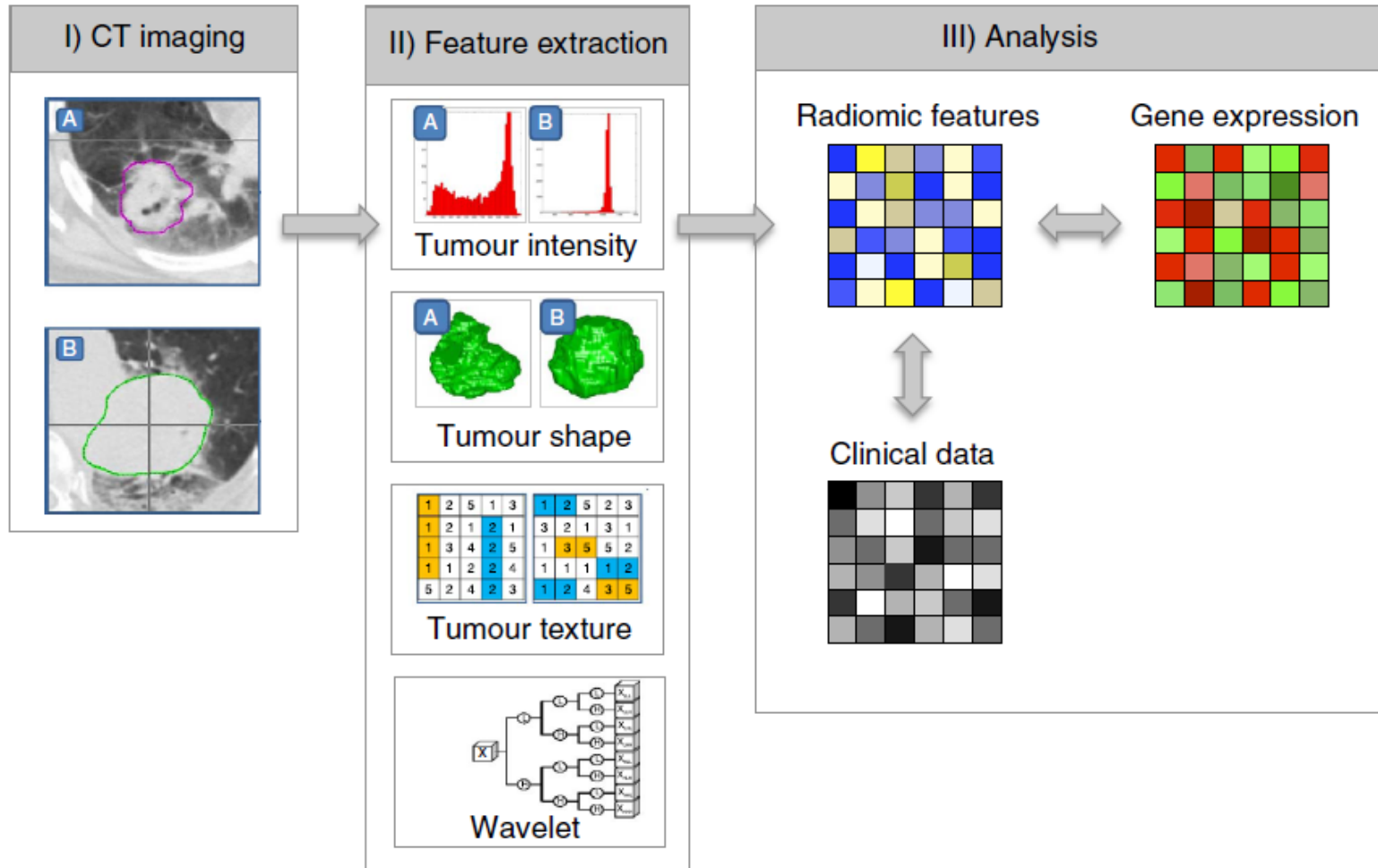
- Characterise
  - Modality
  - Anatomical region
  - Contrasting agent
  - Gender
  - Age
- QC
  - Correct person
  - Missing slices
- Feature detection
  - Artefacts
  - Anomalies

Glocker et al, Vertebrae Localization in Pathological Spine CT via Dense Classification from Sparse Annotations, in MICCAI, September 2013

Criminisi et al, Regression Forests for Efficient Anatomy Detection and Localization in Computed Tomography Scans, in Medical Image Analysis (MedIA), Elsevier, 2013

Criminisi et al, A Discriminative-Generative Model for Detecting Intravenous Contrast in CT Images, in MICCAI, September 2011.

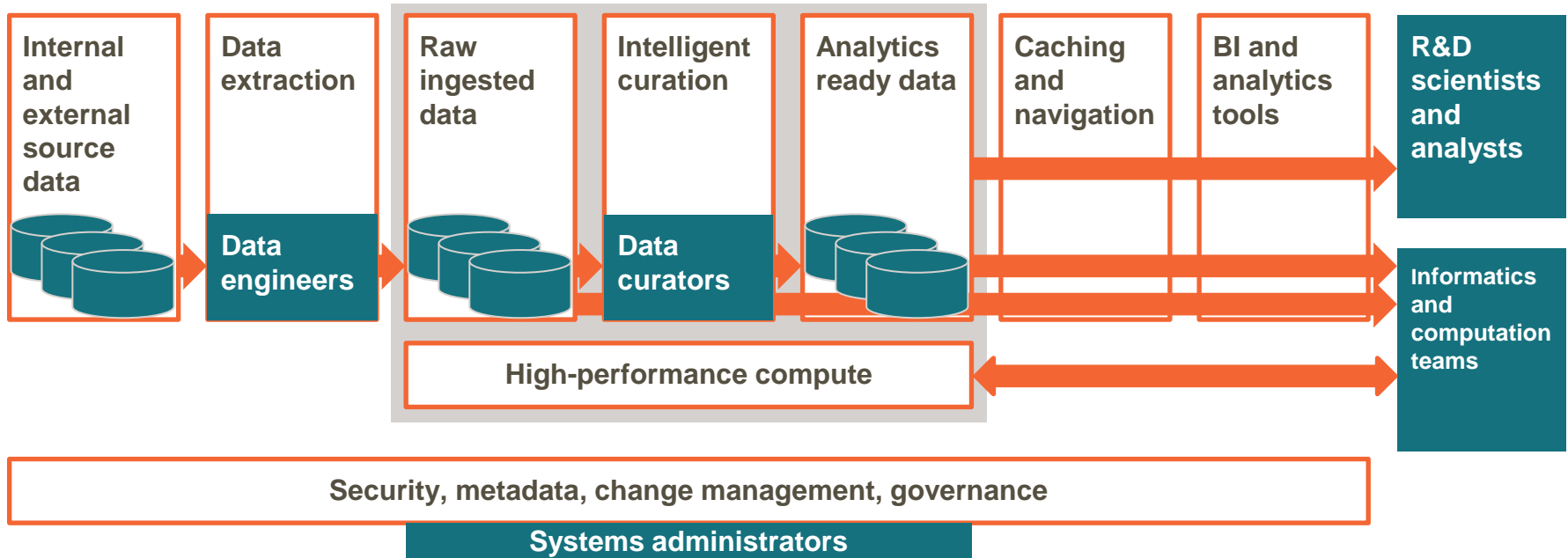
*Detailed quantitative biomarkers are better predictors of survival?*



# Integrative Data Analytics at GSK



*Using technology to make data more accessible*



- A scalable analytics platform for GSK R&D based on Hadoop infrastructure and supporting analytics tools
- Facilitates study of information brought together from multiple domains to uncover unique and actionable insights



- 
- Complex datasets include not only ‘Google style’ big data (i.e. billions of samples) but also other rich datasets (i.e. many data points but not necessarily large numbers of samples)
  - The pharmaceutical industry still relies on very simple analysis methods
  - There is significant scope for improvement!