

# Capturing the Laws of (Data) Nature

Hannes Mühleisen, Martin Kersten & Stefan Manegold

# Statistical Model Fitting & DB?

# Database

# Stats

User gave me a model, let's see.

I am storing some data.

I need some of the observations to fit the model.

This other guy is reading some of my data.

Cool, the model seems to fit the data well!

Let's get some more data to validate the fit...

This other guy is reading some more of my data.

Amazing, model fit is validated.

Beer!

I am storing some data.

# The point?

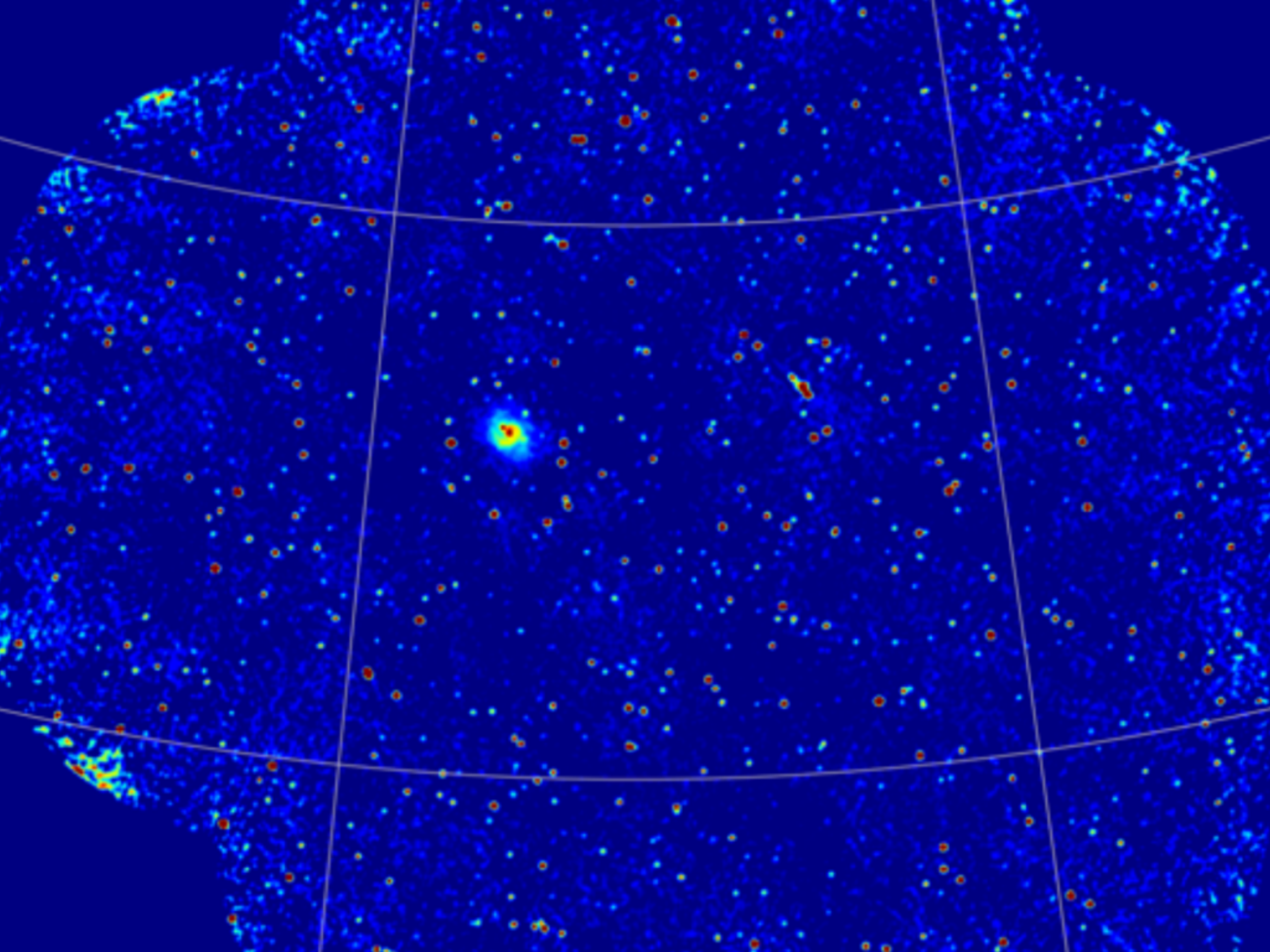
- Everyone has models, **they encode our understanding of the world**
- Everyone has data to train/fit and validate a model
- So far, data management community has **ignored** these models
  - But they hold precious domain knowledge!

# LOFAR Example









Configuration			Measurement	
source	ngz		flux	
1	0.15595		0.23159	
1	0.12392		0.34782	
1	0.14892		0.15927	
1	0.14892		0.32261	
1	0.15595		0.35390	
1	0.14892		0.32784	
1	0.12392		0.33068	
1	0.18486		0.28135	
1	0.18486		0.17872	
1	0.15595		0.28813	

Configuration

Measurement



Grouped by-source operation

|

```
> alpha <- sapply(split(rsmdata, rsmdata$source),
+ function(df) {
+   tryCatch({
+     fit <- nls(flux ~ p * ngz ^ alpha,
+               data=df, start=list(alpha=-.7, p=0.8),
+               control = list(maxiter = 500))
+     coefs <- coef(fit)
+     coefs[[1]]
+   }, error=function(e) {return(as.numeric(NA))
+ })
+ })
```

Model!

Convergence Hints

Measurement

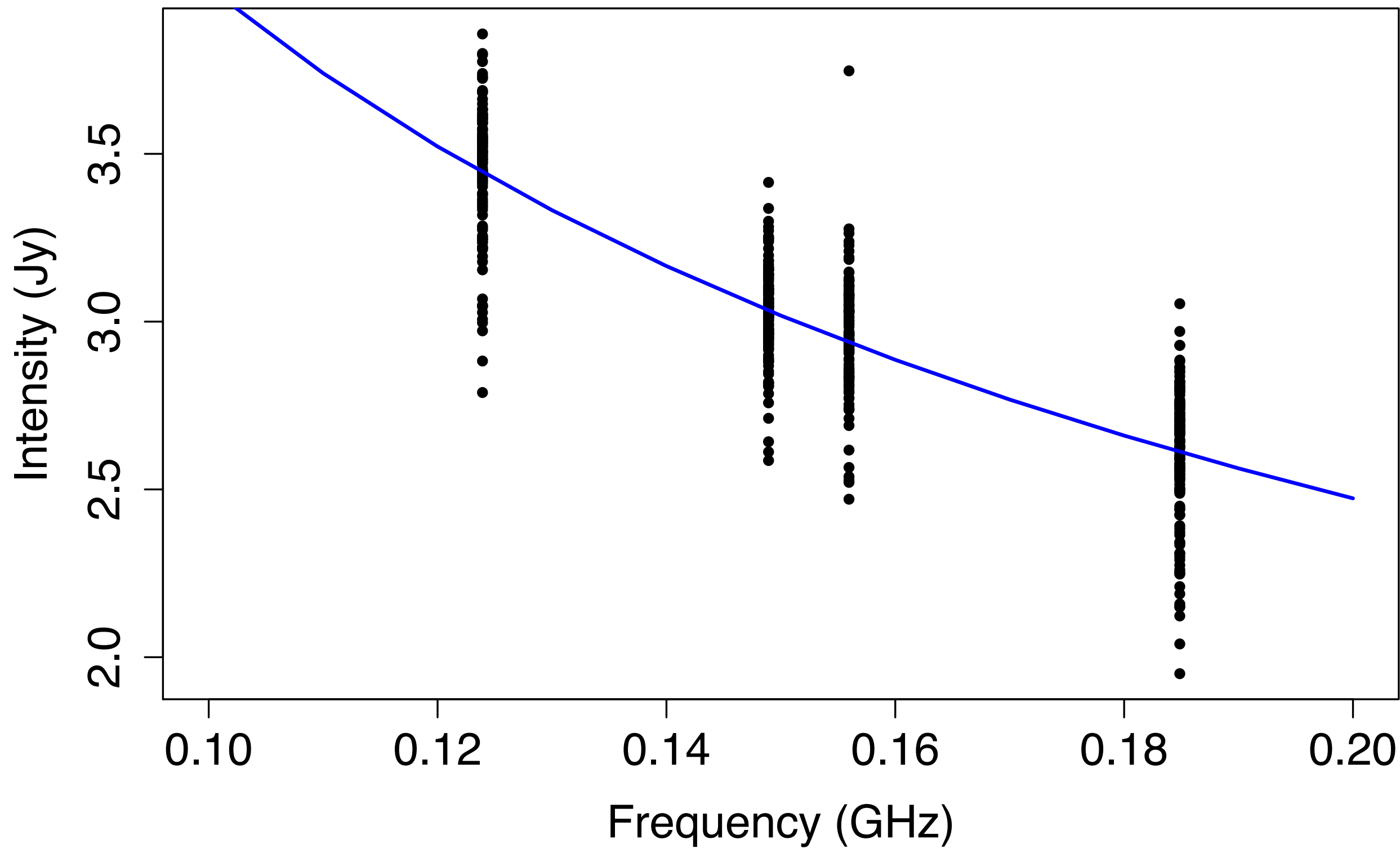
Configuration

`flux ~ p * ngz ^ alpha`

Fitted parameters

Source	Wavelength $\nu$	Intensity $I$		Source	Spectral Index $\alpha$	Constant $p$	Residual SE
1	0.1559555	0.2315911	$\Rightarrow$	1	-0.7183309	0.06257838	0.006559710
1	0.1239243	0.3478159		2	-0.8932245	0.07195620	0.008007786
1	0.1489243	0.1592717		3	-0.7880774	0.56190180	0.016778232
[1,452,821 more rows]				[35,681 more rows]			





Exploit!

```
flux ~ p * ngz ^ alpha
```

↓ Model to function conversion (automatic)

```
flux <- function(p, ngz, alpha) { p * ngz ^ alpha }
```

↓ Move to DB (automatic)

```
CREATE FUNCTION flux3()  
RETURNS TABLE (source integer, nu double, flux DOUBLE) LANGUAGE R {  
  r <- data.frame(...) # model params  
  r$flux <- r$p * r$nu ^ r$alpha  
  return(r[c("source", "nu", "flux")])  
};
```



```
sql>select flux from flux3() where source=17562 and nu=0.156;
```

```
+-----+  
| flux                                     |  
+=====+  
|          2.9370666853479666          |  
+-----+
```

Approximate Answer with zero IO\*

# But...

- What do we do if model parameters are not specified in the query?
  - Sample data?
- Given multiple parameters, it is far from certain that all combinations of values are allowed in the model.
  - Construct filter?

# “Semantic” Compression

	Flux	Flux Residuals	Ratio
ORIG	11,665,408	11,665,408	0%
GZIP	4,331,782	3,748,872	86%
BZIP2	3,341,574	2,752,044	82%
XZ	2,887,584	2,727,144	94%

Drop residuals = lossy compression =





# Data & Model Changes

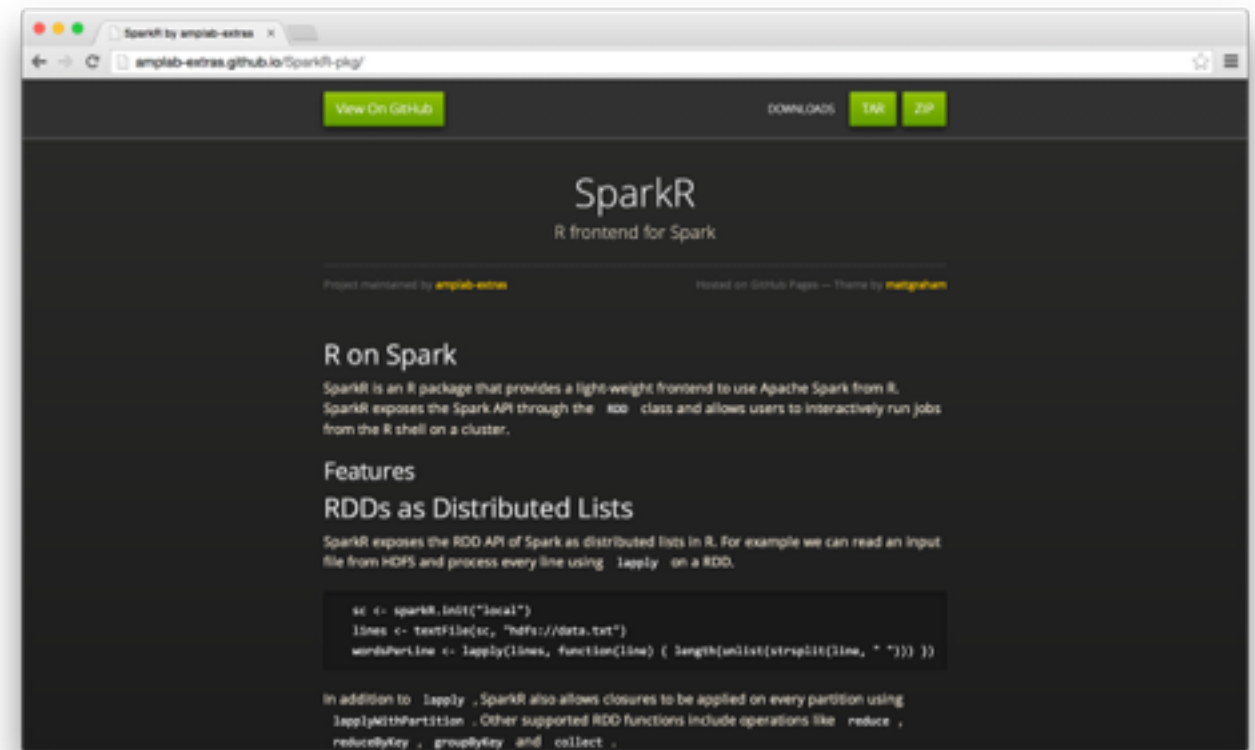
- What should we do if the user gives us a better model?
  - Recompressing could be very expensive
  - Threshold for improvement?
- Changes in the data affect the model quality, too
  - Switch models?
  - Constant Monitoring?

# Multiple, partial or grouped

- There could be many models for a table with overlapping parameters
  - Which one to pick?
- Models do not have to cover the entire table/column
  - “Patching”?
- Models could be fitted on aggregation results
  - Keep group counts?

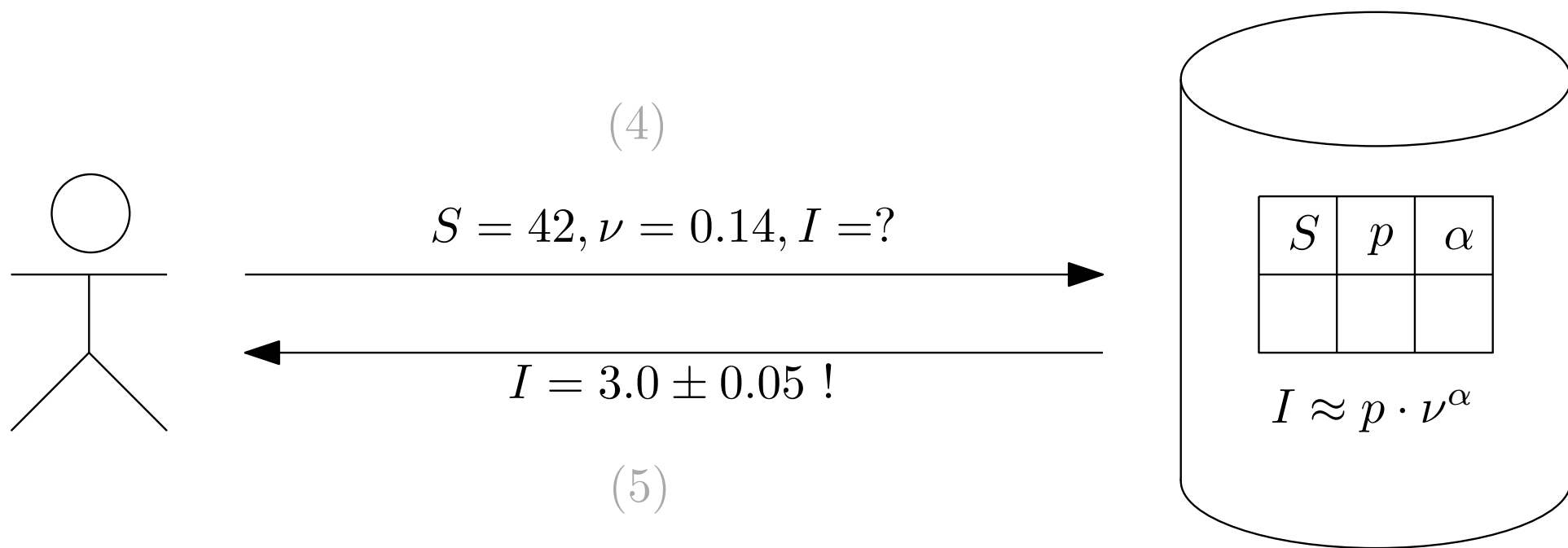
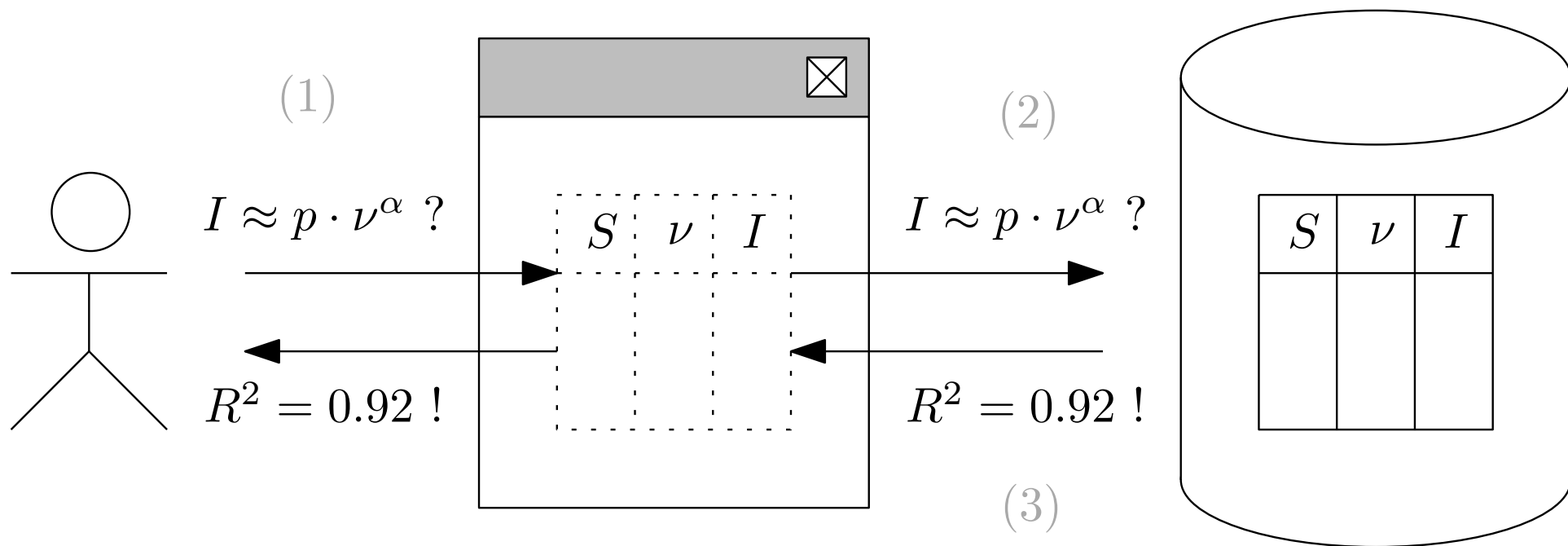
How do we get our  
hands on Models?





# Integrate & Intercept

- Integrate model fitting infrastructure into data management system.
  - Also: **Huge** performance benefits for analysts!
- Intercept model fitting and validation operations by the user and store the model for later use.
  - Storage format: Model code + Parameters



“Essentially, all models are wrong,  
but some are useful.”

George E. P. Box

Questions?