

**ETH**



# **Standardized Tests of RELM ERF's Against Observed Seismicity**

**Danijel Schorlemmer<sup>1</sup>, Dave D. Jackson<sup>2</sup>,  
Matt Gerstenberger<sup>3</sup>, Stefan Wiemer<sup>1</sup>**

<sup>1</sup> ETH Zürich, Swiss Seismological Service, Switzerland

<sup>2</sup> Department of Earth and Space Sciences, UCLA, USA

<sup>3</sup> Institute of Geological & Nuclear Sciences, New Zealand

## Why do we test?

- Agreement of a group
- Unbiased forward-looking test
- Fully specified in advance
- Rigorous test

## What do we test?

- Comparative performance of ERF's
- Consistency of ERF's with observations

## How do we test?

## Examples from Parkfield

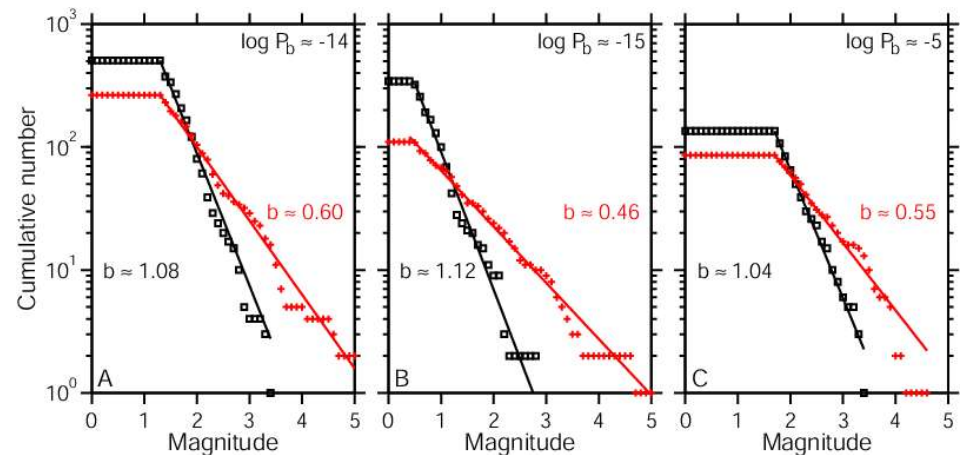
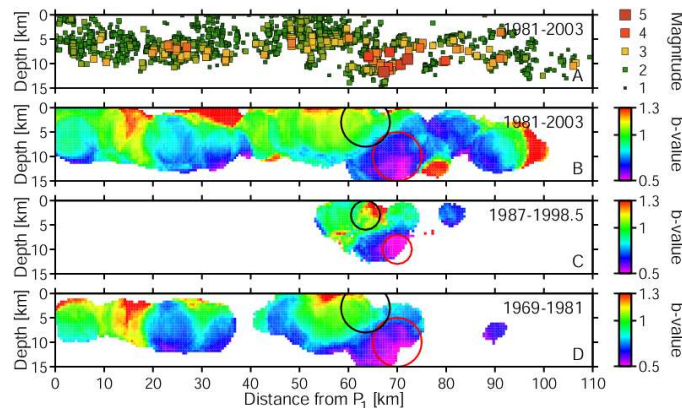
# Example – Test

## Model 1 (Variable b-value)

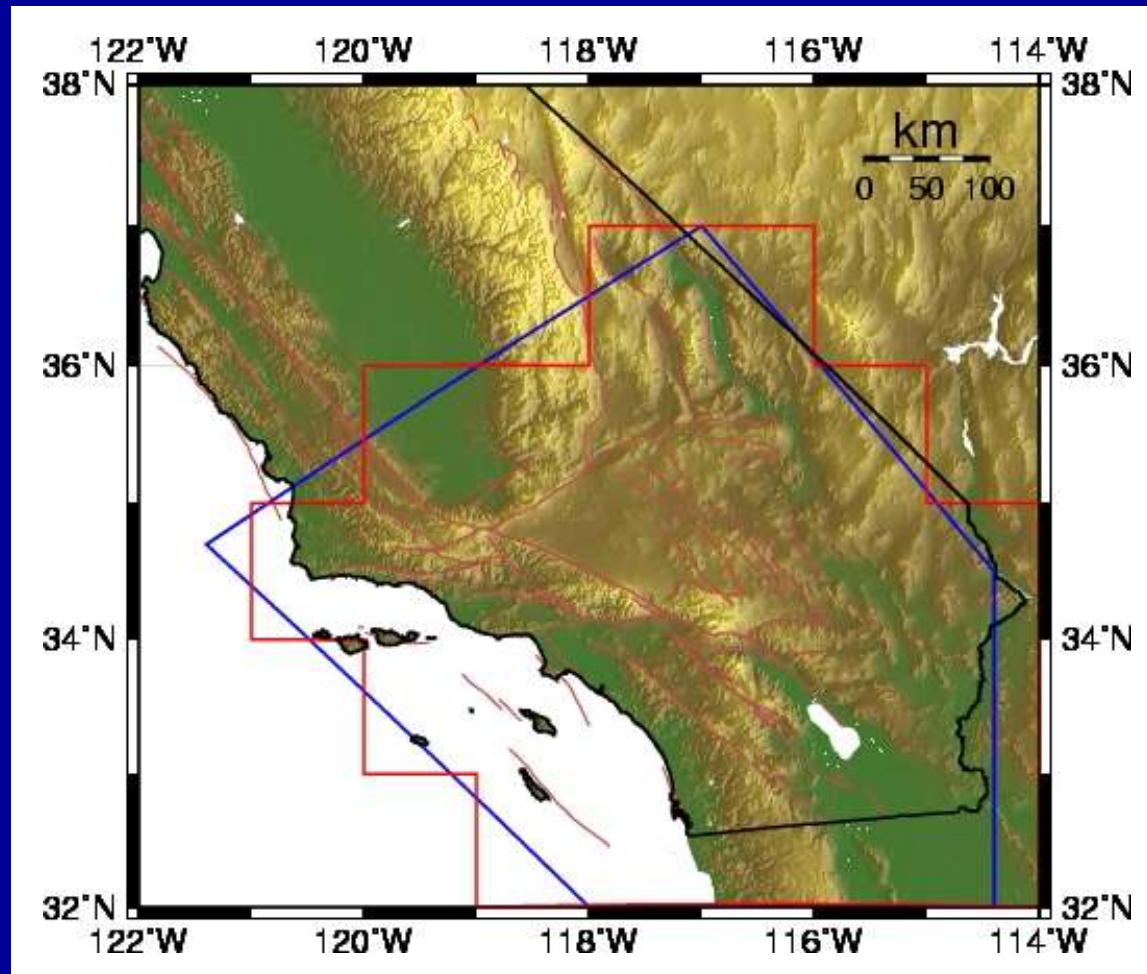
Forecast generated by extrapolating the frequency-magnitude distribution using spatially varying a- and b-values

## Model 2 (Constant b-value)

Same as model 1, but with regional constant b-value



## Testing area in California

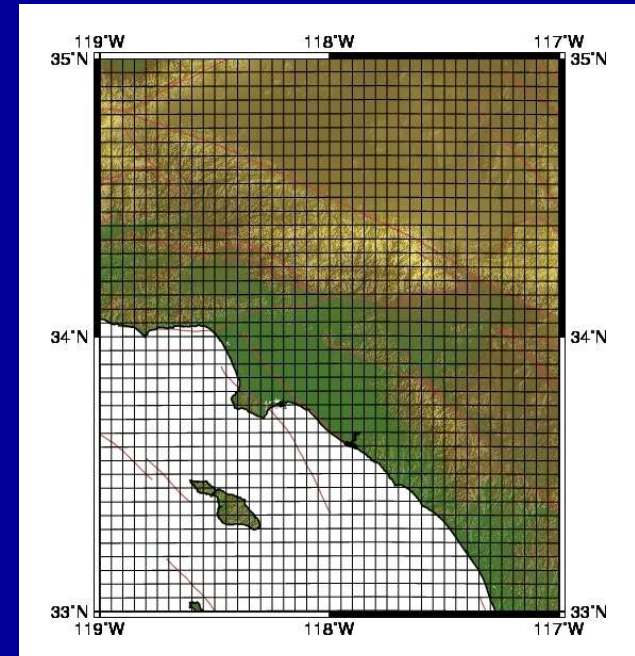


What is a bin?

A bin  $b_i$  defines a volume, magnitude range, period, and range of focal mechanism angles for which a forecast is issued

Our proposed default binning:

Lon/Lat	$0.05^\circ \times 0.05^\circ$
Depth	no binning: 0-30km
Magnitude	0.1
Focal M.	$30^\circ$



A forecast of a model  $j$  is defined as the expectations  $\lambda_i^j$  per bin  $b_i$

$$\forall b_i | \lambda_i^j := \lambda_i^j(b_i)$$

The total expectation of a model  $j$  is the vector  $\Lambda^j$

$$\Lambda^j = \begin{pmatrix} \lambda_1^j \\ \lambda_2^j \\ \vdots \\ \lambda_n^j \end{pmatrix}$$

Forecasts are issued for

**1 year**      for quasi-stationary models

**1 day**        for time-dependent models

Forecasters provide expectations, i.e. numbers of events per bin, not programs

# Observations

Observations  $\omega_i$  (number of earthquakes) per bin  $b_i$  are a vector  $\Omega$

$$\Omega = \begin{pmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_n \end{pmatrix}$$

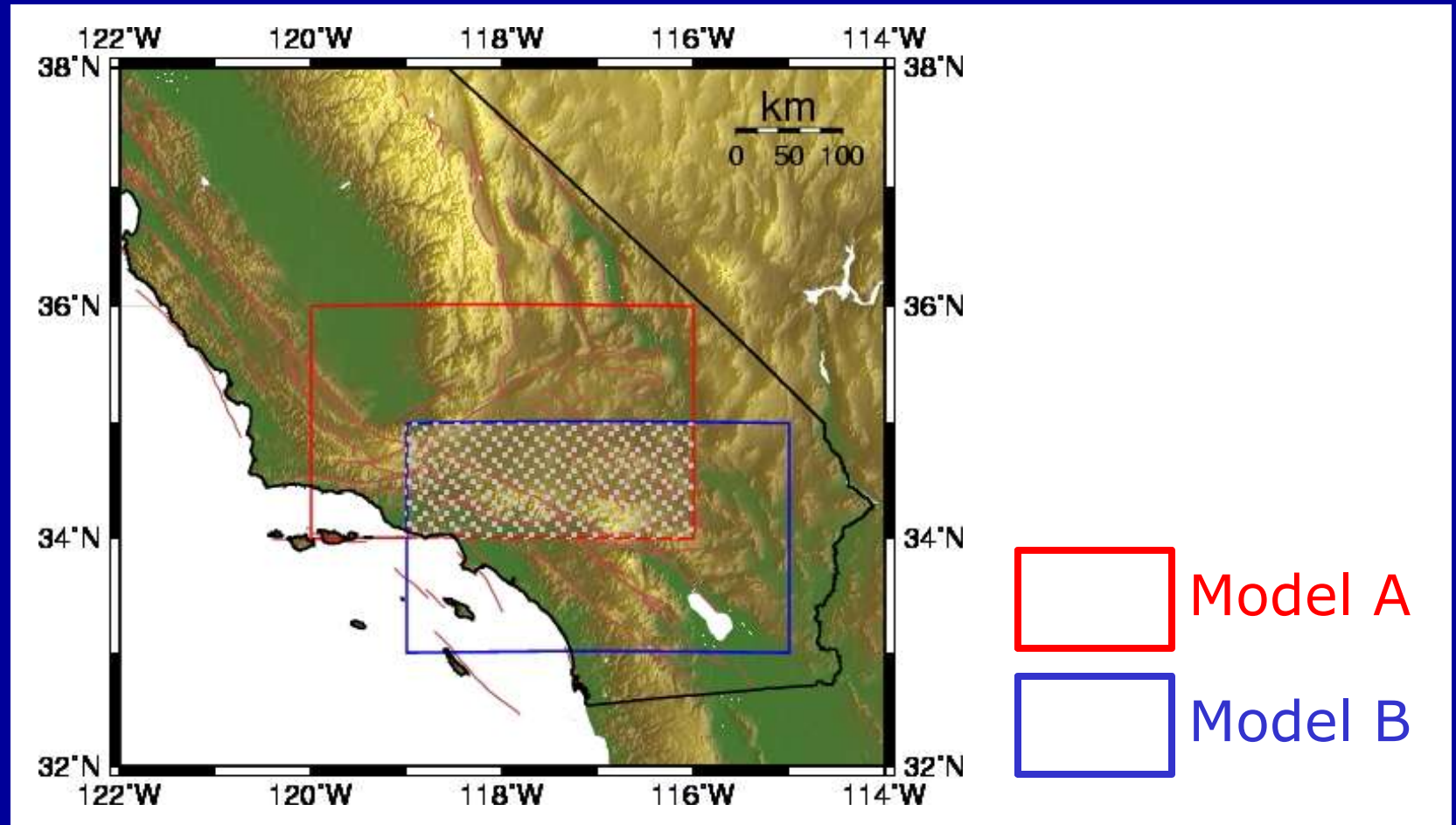
We use 2 catalogs:

1. Y. Kagan's RELM catalog for hypocenter, magnitude, focal time
2. Harvard catalog for focal mechanisms



# Limitations

We can only test in corresponding bins of two models





# Agreements

ETH

Forecast is a vector of expectations (rates): quakes per year (or day) per bin

Forecasters provide expectations, not programs

All quakes count: no distinction between foreshocks, main shocks, and aftershocks

Bins of  $0.05^\circ \times 0.05^\circ \times 0.1$  magnitude units

Two main “menu items”:

- Five year forecast of  $M \geq 5$ , no updates
- Five year forecast of  $M \geq 4$ , updated daily

Special orders ok if

- Multiple models use same bins, and
- Sufficient earthquakes for test

Likelihood is based on the poissonian probability distribution

$$p(\omega|\lambda) = \frac{\lambda^\omega}{\omega!} e^{-\lambda}$$

Log-Likelihood of model  $j$  in bin  $b_i$  of the observation  $\omega_i$  given the expectation  $\lambda_i^j$

$$L(\omega_i|\lambda_i^j) = -\lambda_i^j + \omega_i \log \lambda_i^j - \log \omega_i!$$

The joint log-likelihood  $L^j$  of model  $j$

$$L^j = L(\Omega | \Lambda^j) = \sum_{i=0}^n L(\omega_i | \lambda_i^j)$$

Calculation example

Forecast bin			Expectation		Probabilities					
Bin	Region	Mag.	Ann. rate	$\lambda$	0	1	2	3	4	etc.
1	1	4.5-5.5	0.40	2.0	0.14	0.27	0.27	<b>0.18</b>	0.09	
2	1	5.5-6.5	0.04	0.2	<b>0.82</b>	0.16	0.02	0.00	0.00	
3	2	4.5-5.5	0.20	1	0.37	<b>0.37</b>	0.18	0.06	0.02	
4	2	5.5-6.5	0.02	0.1	<b>0.90</b>	0.09	0.00	0.00	0.00	

Likelihood  $L = \log(0.18 * 0.82 * 0.37 * 0.90)$

# Testing – Likelihood-ratio

The log-likelihood-ratio  $R$  of two models

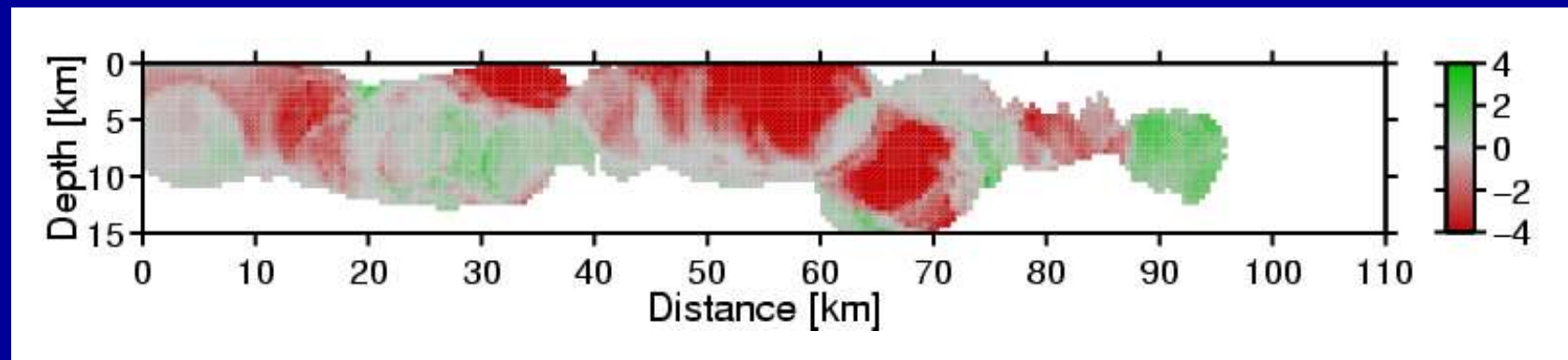
$$R = L^2 - L^1 = L(\Omega|\Lambda^2) - L(\Omega|\Lambda^1)$$

First evaluation of results

$R < 0$  Higher likelihood of **model 1 (Variable  $b$ )**

$R > 0$  Higher likelihood of **model 2 (Constant  $b$ )**

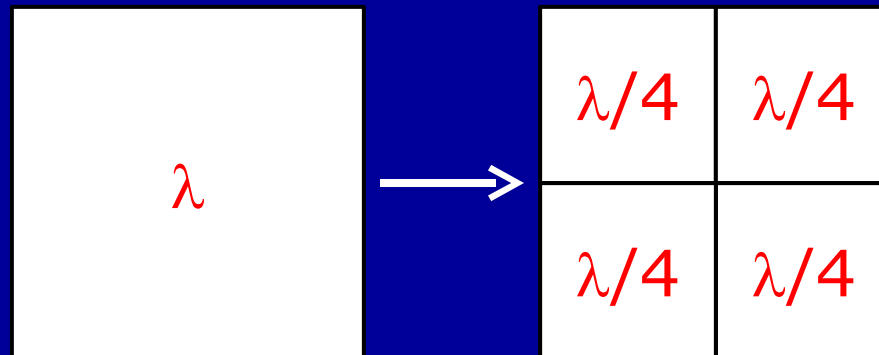
Likelihood-ratio per bin at Parkfield



# Testing – Different Resolutions

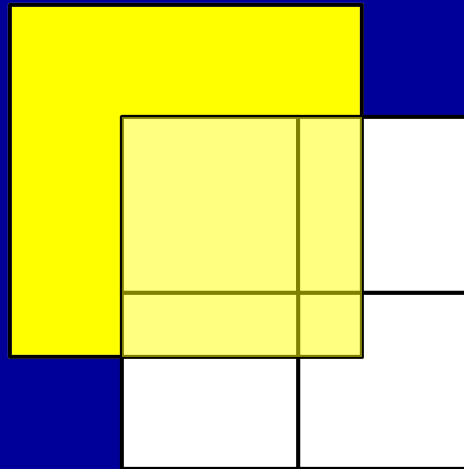
Testing of models with different resolutions of bins

Distribute the expectations of the low resolution forecast over the bins of the high resolution forecast



This approach does not change the likelihood-ratio

## Constraints of this approach



Limits of bins need to be aligned

## **N-Test**

Test if the number of observed events is in the range of the expectation of a model

## **L-Test**

Examines the consistency of a model with the observation

## **R-Test**

Compares 2 models by its log-likelihood-ratio





ETH

# Testing – Simulations

Create 'simulated catalogs'

for every bin draw a random number from uniform distribution  $[0,1]$

create simulated number of events based on random number and expectation of the bin (Poissonian distribution)

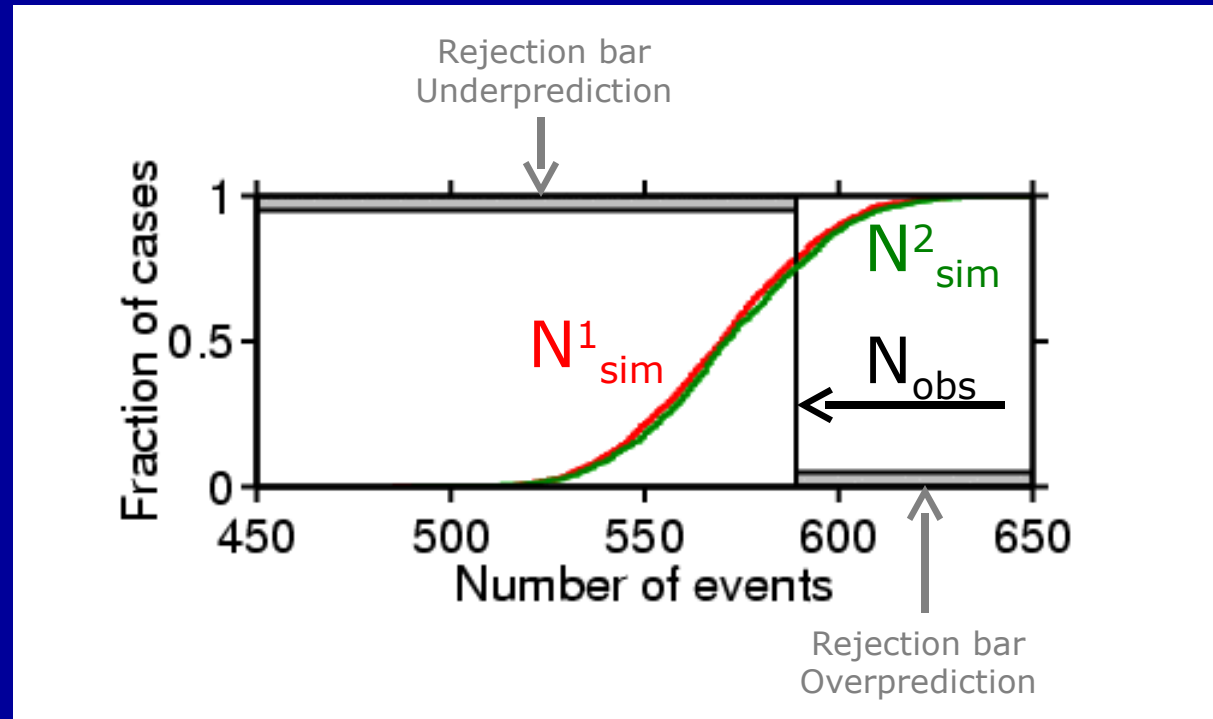
repeat over all bins

repeat 1000 times for the set of bins

→ Simulation based on expectations of a model

# N-Test

Compare total number of observed events  $N_{obs}$  with number of simulated events  $N_{sim}^j$



If the distribution  $N_{sim}^j$  touches one of the rejection bars, the model is over- or underpredicting events

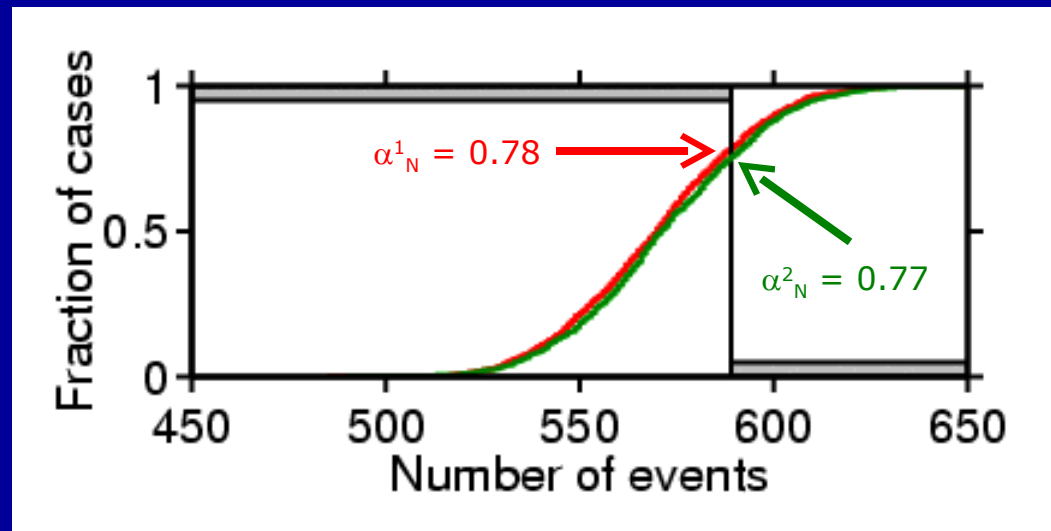
# N-Test

Quantify performance of a model with  $\alpha_N^j$ , fraction of  $N_{sim}^j < N_{obs}$

Reject model if

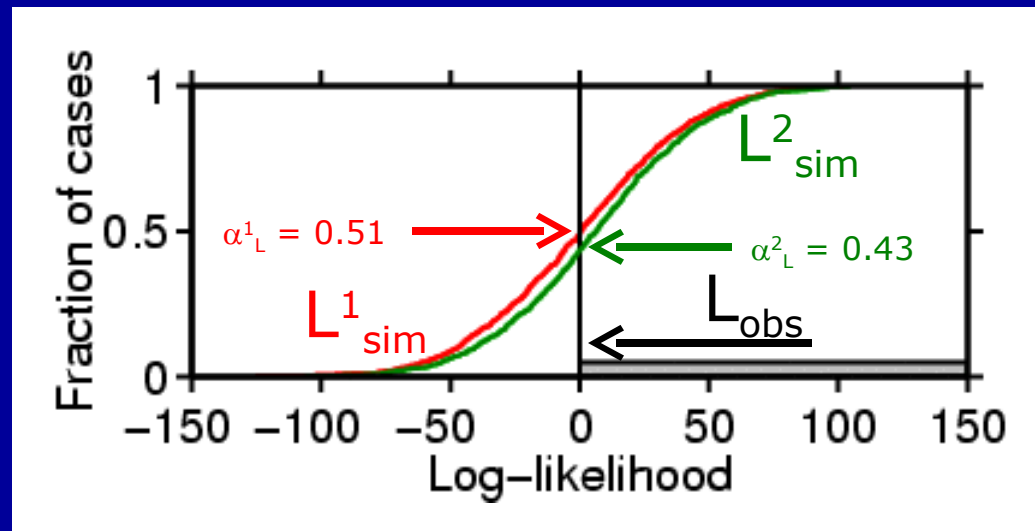
$$\alpha_N^j > 0.95 \text{ or}$$

$$\alpha_N^j < 0.05$$



This test shows whether the observed number of events is in the range of the model's forecast

Compare the likelihood  $L_{obs}^j$  of a model with the likelihoods  $L_{sim}^j$  of simulated catalogs



We normalize the log-likelihoods  $\rightarrow L_{obs}^j = 0$

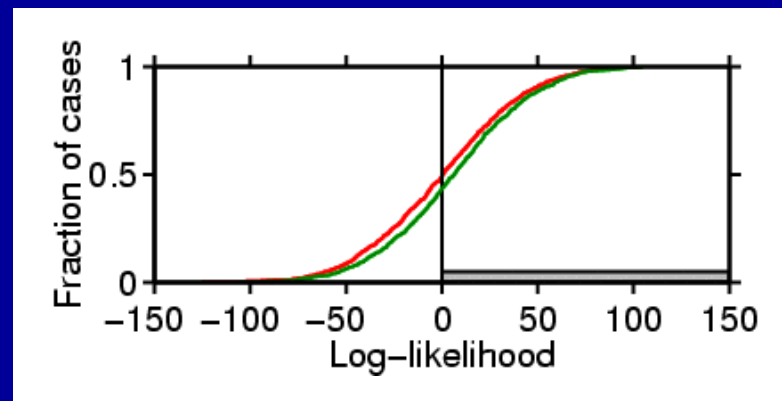
This test shows whether the model's forecast is consistent with the observation

Why can't we reject  $\alpha_L^j > 0.95$ ?

Underprediction vs. 'perfect forecast'

A 'perfect forecast' would have a higher likelihood than any of the simulated likelihoods -> We cannot reject a 'perfect forecast'

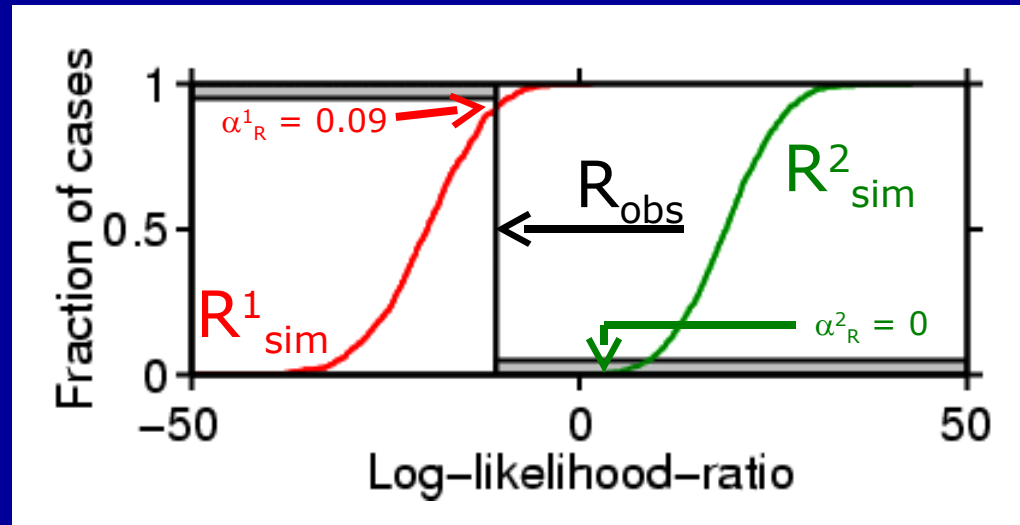
Models underpredicting events may succeed in the R-Test → N-Test



Compare the likelihood-ratio  $R_{obs}$  of two models with the likelihood-ratios  $R_{sim}^1$  and  $R_{sim}^2$  of simulated catalogs based on the expectations of both models

$R_{sim}^1$  simulation based on **model 1**

$R_{sim}^2$  simulation based on **model 2**



# Evaluation

Evaluate the likelihood-ratio of pair-wise tests of all models  
→ R-Test

Test for consistency of the models with the observation  
→ L-Test, N-Test

Analyze the spatial performance of models

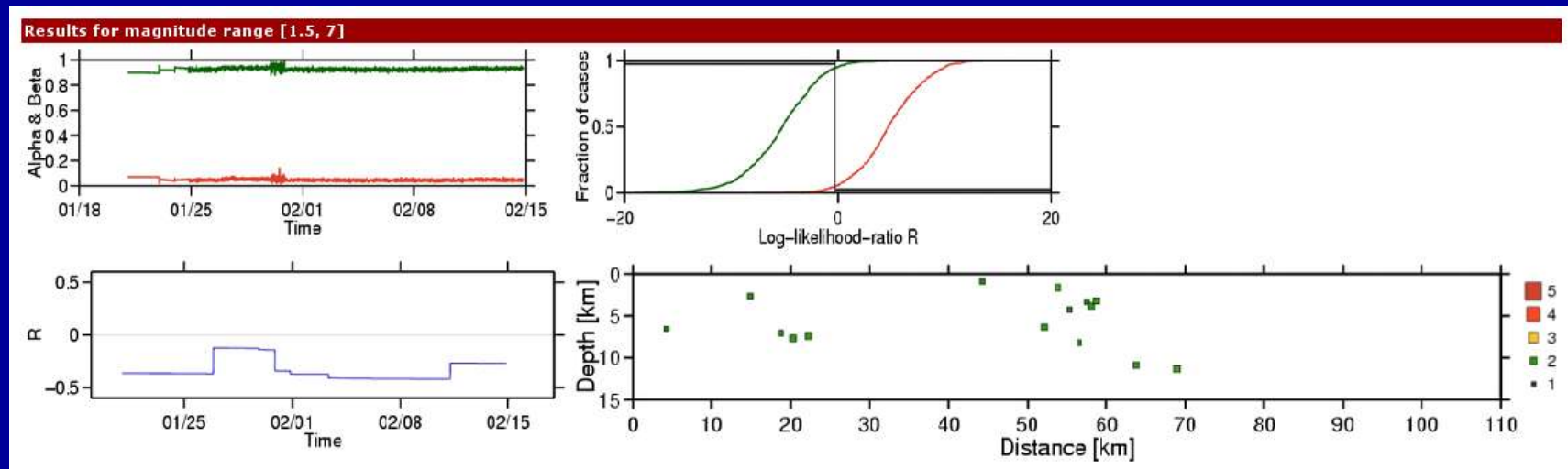
# Conducting the Tests

Compute the tests daily to track the temporal performance

Evaluate yearly to determine the overall and spatial performance

Yearly model investigations

Parkfield realtime test





Is the Poissonian distribution appropriate for likelihood computation?

→ Negative binomial distribution

Should we use declustered catalogs for quasi-stationary models?



Thank you

ETH

Track the Parkfield realtime test from next month  
on at:

<http://cool.ethz.ch>