

# **Data Analysts Without Borders: A Case Study**

Warren G. Anderson  
University of Wisconsin - Milwaukee



Image courtesy of NASA  
Goddard NR group

# Outline

- Foundation
  - Four Data Analysis (DA) groups
- Formulation
  - GWs, IFOs, Data, Likelihood, Power, Power<sup>2.0</sup>
- Application
  - Searches for modeled signals, backgrounds, bursts
- Migration
  - Real life, burst or glitch, STAMP
- Conclusion

# Foundation

- Ground-based gravitational wave DA is traditionally divided into four categories:

Category		

# Foundation

- Ground-based gravitational wave DA is traditionally divided into four categories:

Category	Short Duration	

# Foundation

- Ground-based gravitational wave DA is traditionally divided into four categories:

Category	Short Duration	Long Duration

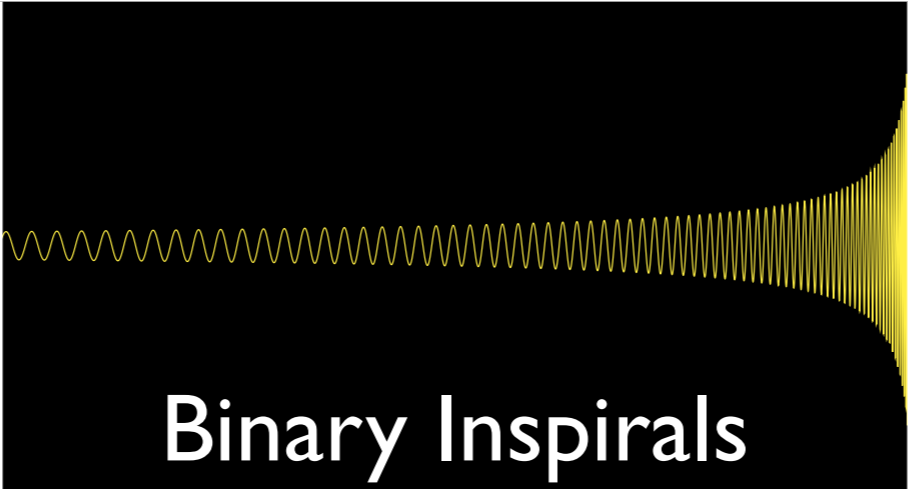
# Foundation

- Ground-based gravitational wave DA is traditionally divided into four categories:

Category	Short Duration	Long Duration
Theoretical Waveform		

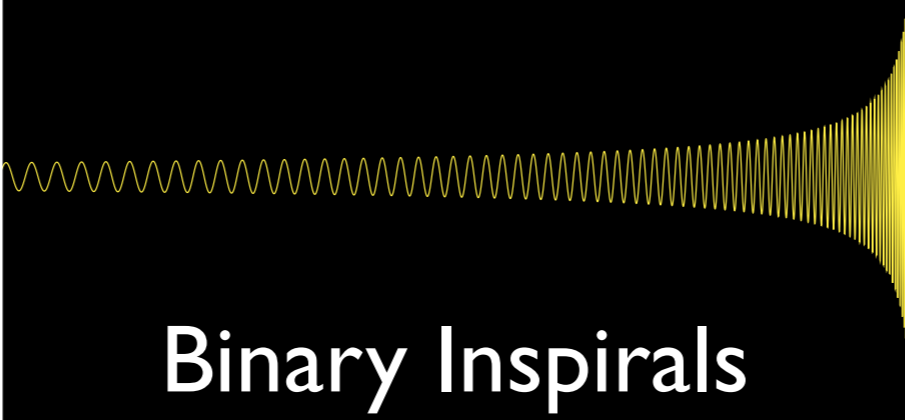
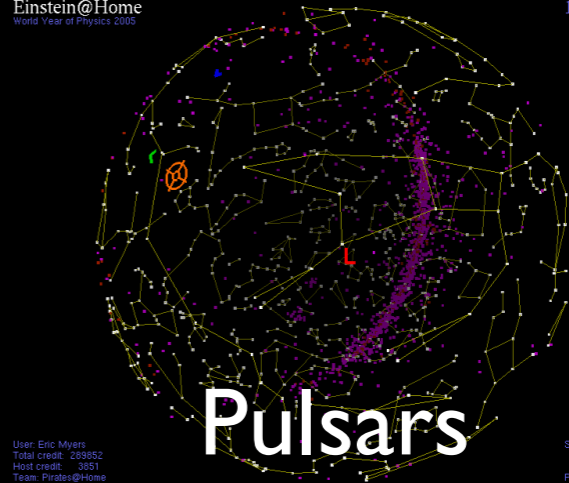
# Foundation

- Ground-based gravitational wave DA is traditionally divided into four categories:

Category	Short Duration	Long Duration
Theoretical Waveform	 <p>Binary Inspirals</p>	

# Foundation

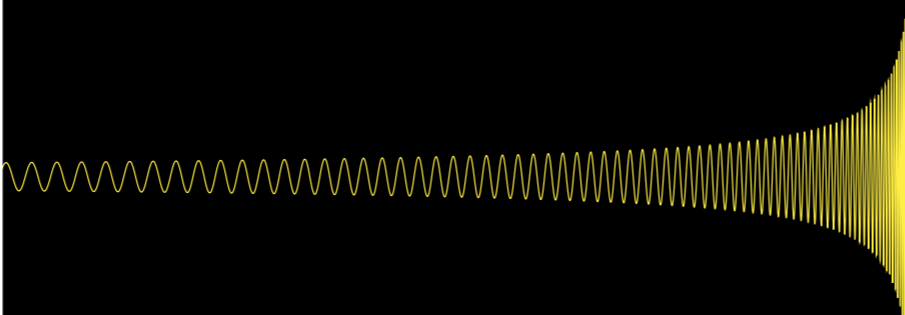
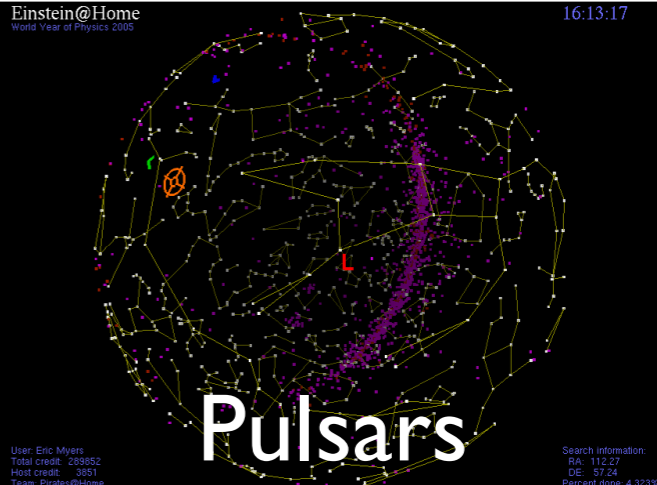
- Ground-based gravitational wave DA is traditionally divided into four categories:

Category	Short Duration	Long Duration
Theoretical Waveform	 <p data-bbox="900 1270 1495 1360">Binary Inspirals</p>	 <p data-bbox="1970 1270 2252 1346">Pulsars</p> <p data-bbox="1778 890 2414 1371"><small>Einstein@Home World Year of Physics 2005 16:13:17 User: Eric Myers Total credit: 26952 Host credit: 351 Team: Pulsar@Home Search information: RA: 112.27 DE: 57.24 Percent done: 4.323%</small></p>



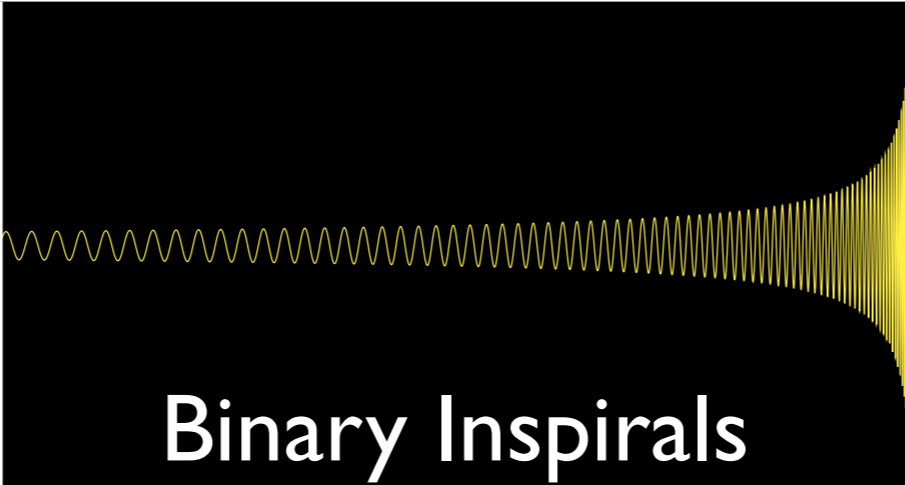
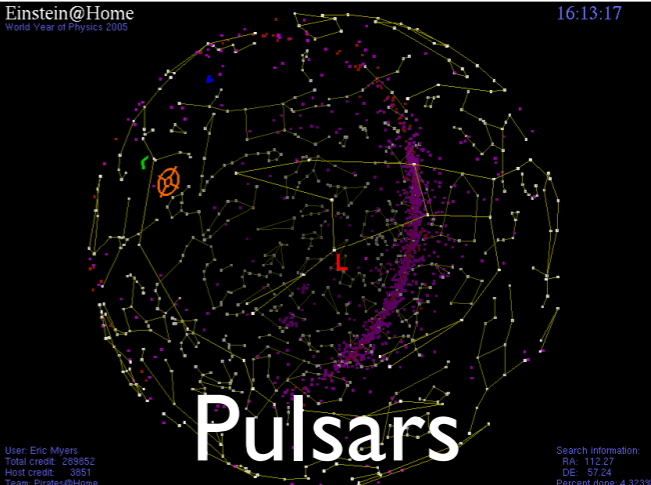
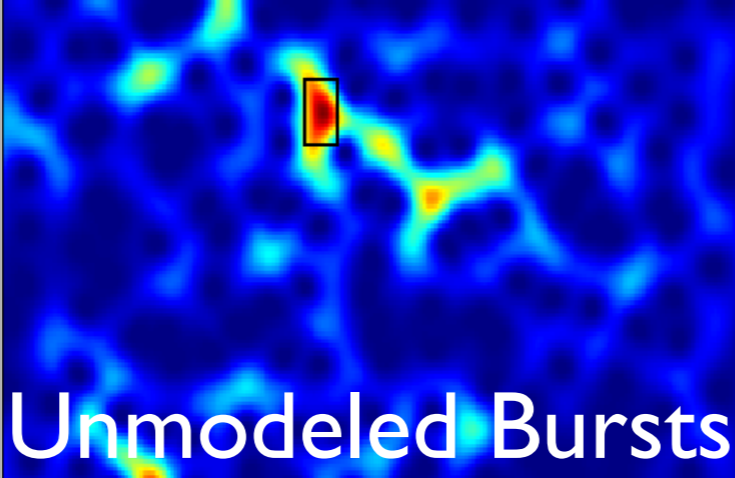
# Foundation

- Ground-based gravitational wave DA is traditionally divided into four categories:

Category	Short Duration	Long Duration
Theoretical Waveform	 <p data-bbox="900 1268 1495 1360">Binary Inspirals</p>	 <p data-bbox="1972 1268 2247 1346">Pulsars</p>
No Theoretical Waveform		

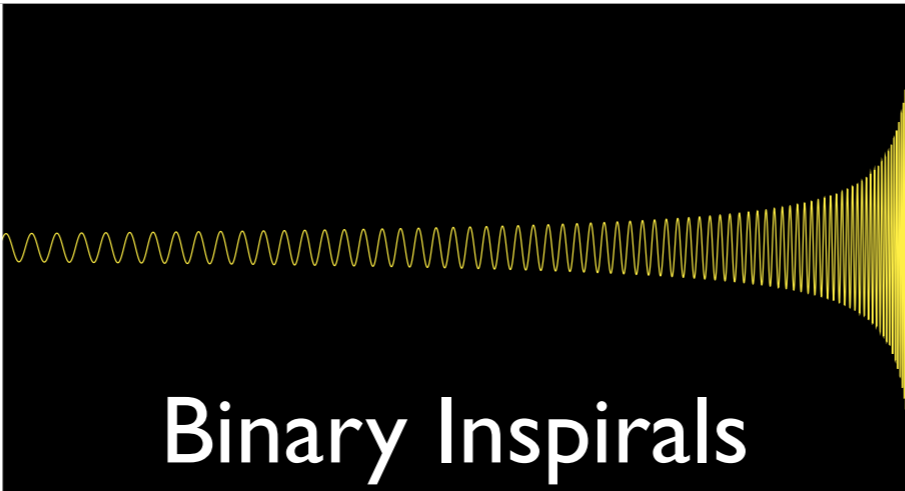
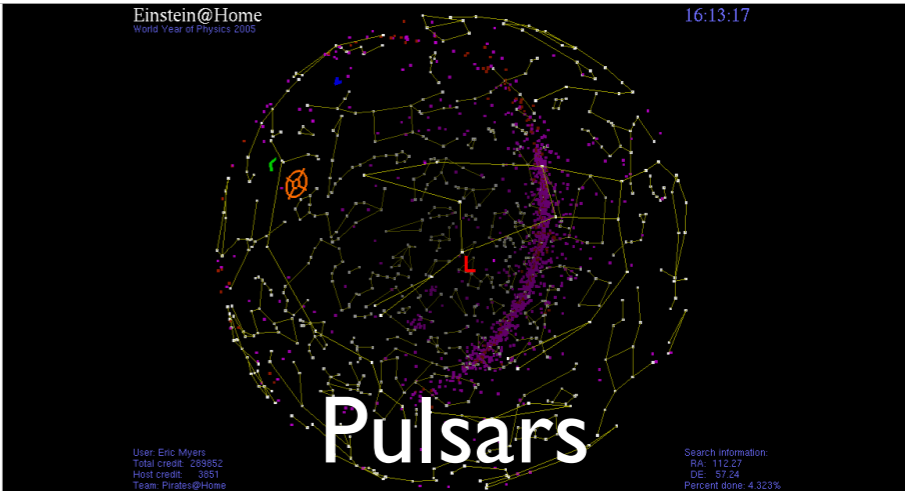
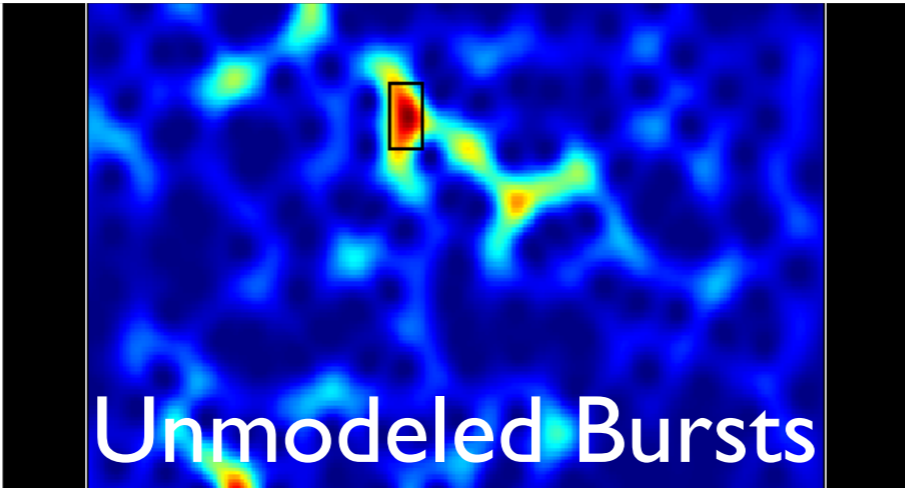
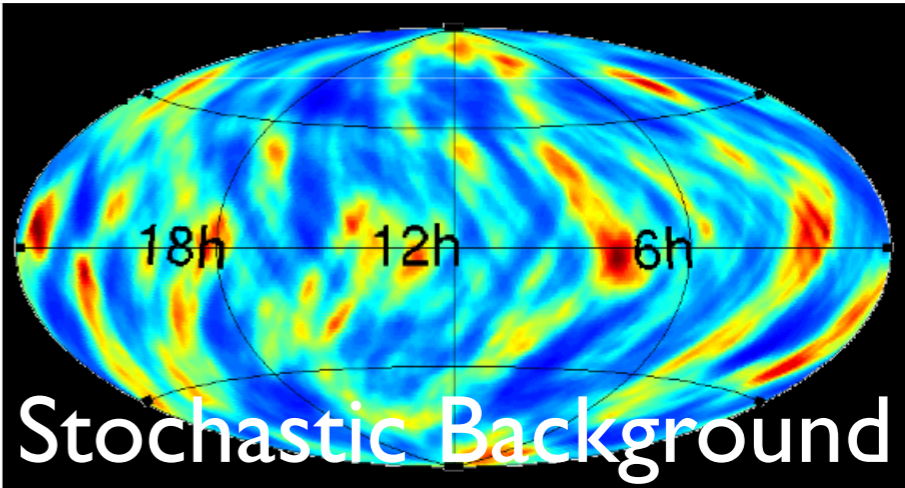
# Foundation

- Ground-based gravitational wave DA is traditionally divided into four categories:

Category	Short Duration	Long Duration
Theoretical Waveform	 <p data-bbox="900 1268 1495 1360">Binary Inspirals</p>	 <p data-bbox="1970 1268 2252 1346">Pulsars</p>
No Theoretical Waveform	 <p data-bbox="834 1759 1561 1837">Unmodeled Bursts</p>	

# Foundation

- Ground-based gravitational wave DA is traditionally divided into four categories:

Category	Short Duration	Long Duration
Theoretical Waveform	 <p data-bbox="900 1270 1495 1360">Binary Inspirals</p>	 <p data-bbox="1970 1270 2244 1348">Pulsars</p>
No Theoretical Waveform	 <p data-bbox="831 1764 1566 1841">Unmodeled Bursts</p>	 <p data-bbox="1668 1764 2546 1841">Stochastic Background</p>

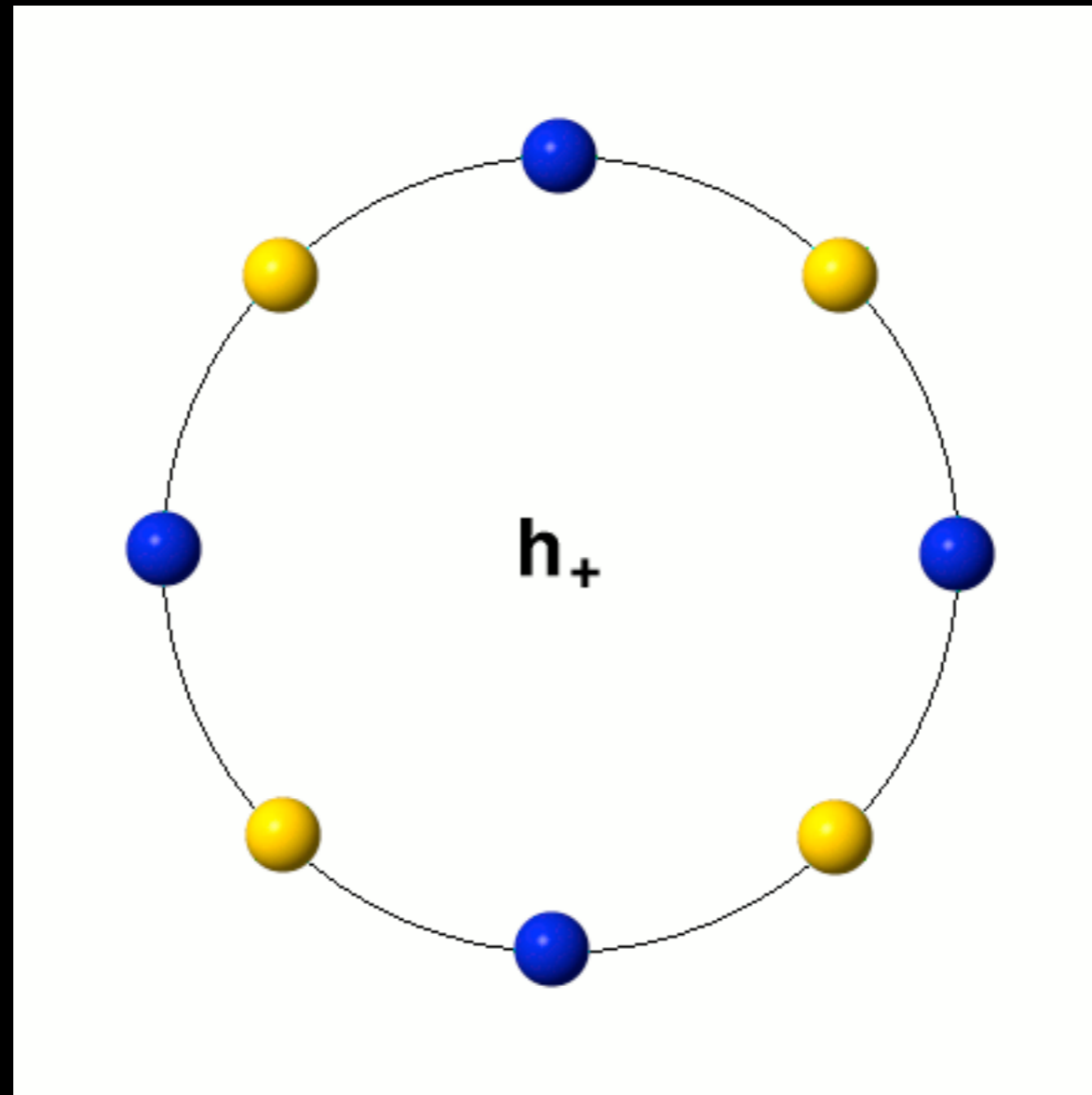
# Foundation

- Data analysts are usually divided into four groups based on the categories of signal they are looking for.
- Questions:
  - “Are these categories fundamental or just convenient?”
  - “Do all signals fit neatly into just one category?”
  - “Is there one search method that is optimal for all signals in each category?”

# Formulation: GWs



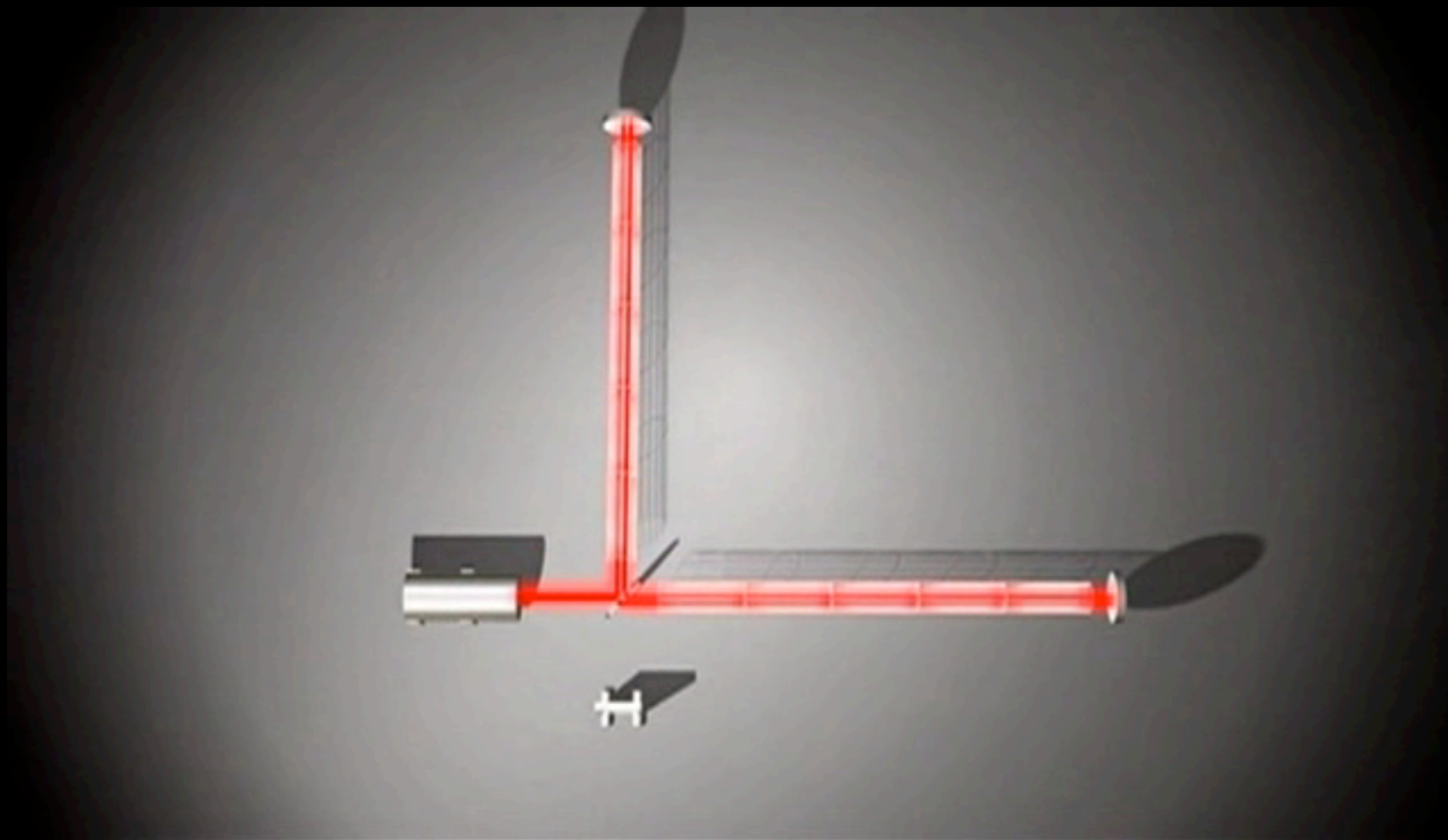
# Formulation: GWs



# Formulation: IFOs

From “Einstein’s Messengers”  
National Science Foundation

# Formulation: IFOs



From "Einstein's Messengers"  
National Science Foundation



# Formulation: Data

- From our instruments we get strain data  $h(t)$

$$h(t) = \underset{\text{Noise}}{n(t)} + \underset{\text{Signal}}{s(t)}$$

# Formulation: Data

- From our instruments we get strain data  $h(t)$

$$h_i = \underbrace{n_i}_{\text{Noise}} + \underbrace{s_i}_{\text{Signal}}$$

- Signal (usually) is deterministic (possibly  $s_i \approx 0$ )
- Noise is stochastic - eg Gaussian (ideally)

$$p(\mathbf{n}) \propto \exp\left(-\frac{\mathbf{n} \cdot \mathbf{n}}{2\sigma}\right)$$

# Formulation: Likelihood

- Neyman-Pearson: Optimal statistic is **Likelihood**

$$\Lambda[h] = \int \frac{p(h|\mathbf{n} + \mathbf{s})}{p(h|\mathbf{n})} \mathcal{D}[\mathbf{s}]$$

- The measure  $\mathcal{D}[\mathbf{s}]$  projects the integrand from the space of all possible signals down onto the subspace of signals we are searching for, denoted by  $\Sigma$ .

# Formulation: Likelihood

- Neyman-Pearson: Optimal statistic is **Likelihood**

$$\Lambda[h] = \int \frac{p(h|\mathbf{n} + \mathbf{s})}{p(h|\mathbf{n})} \mathcal{D}[\mathbf{s}]$$

*Signal*

- The measure  $\mathcal{D}[\mathbf{s}]$  projects the integrand from the space of all possible signals down onto the subspace of signals we are searching for, denoted by  $\Sigma$ .

# Formulation: Likelihood

- Neyman-Pearson: Optimal statistic is **Likelihood**

$$\Lambda[h] = \int \frac{p(h|\mathbf{n} + \mathbf{s})}{p(h|\mathbf{n})} \mathcal{D}[\mathbf{s}]$$

*Signal*

- Probabilities of  $h$  derive from distribution of  $\mathbf{n}$ , eg

$$p(h|\mathbf{n} + \mathbf{s}) \propto e^{-\frac{(\mathbf{h} - \mathbf{s}) \cdot (\mathbf{h} - \mathbf{s})}{2\sigma}}$$

# Formulation: Likelihood

- Neyman-Pearson: Optimal statistic is **Likelihood**

$$\Lambda[h] = \int \frac{p(h|\mathbf{n} + \mathbf{s})}{p(h|\mathbf{n})} \mathcal{D}[\mathbf{s}]$$

*Noise*                      *Signal*

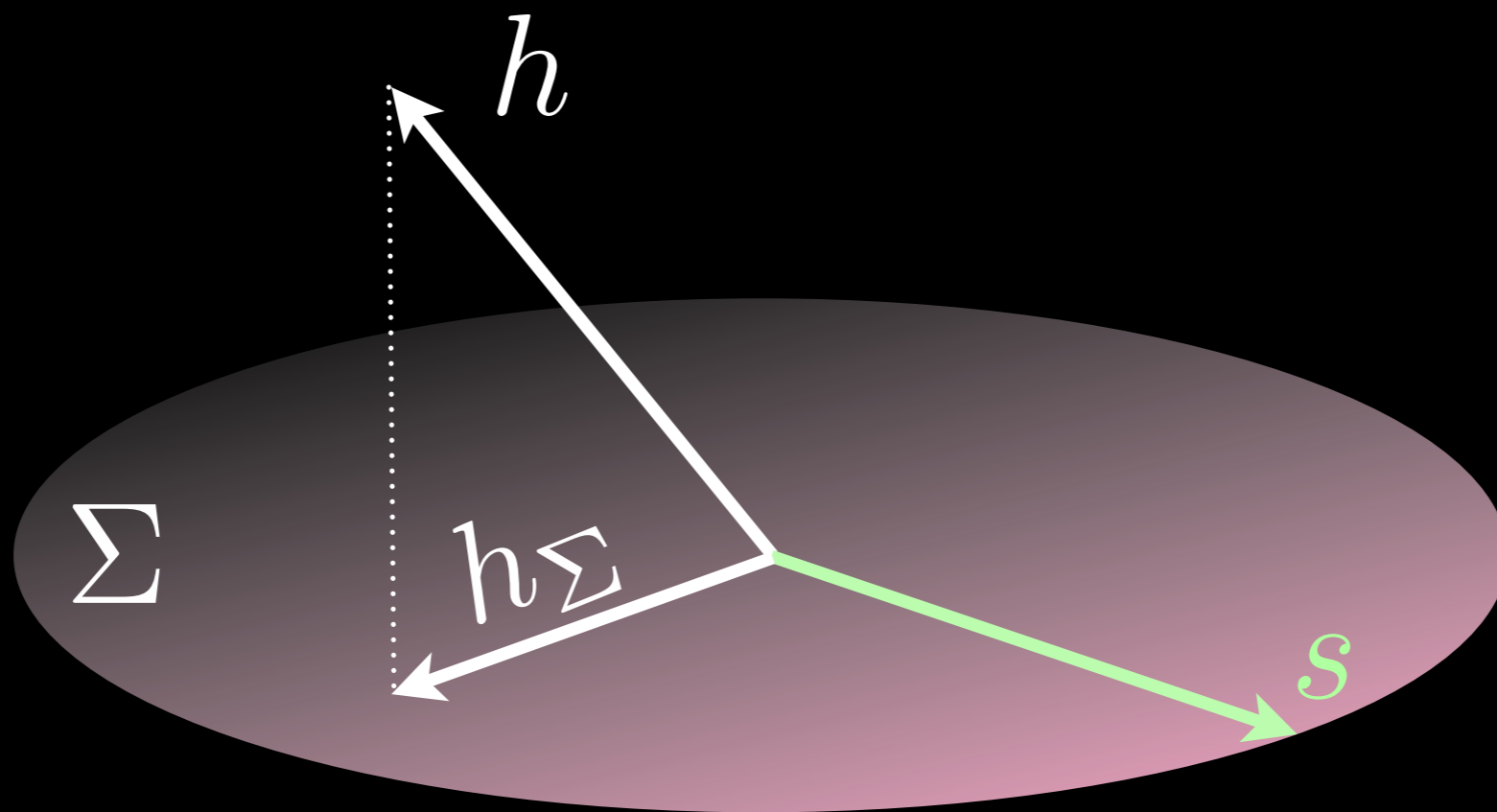
- Probabilities of  $h$  derive from distribution of  $\mathbf{n}$ , eg

$$p(h|\mathbf{n} + \mathbf{s}) \propto e^{-\frac{(\mathbf{h} - \mathbf{s}) \cdot (\mathbf{h} - \mathbf{s})}{2\sigma}}$$

# Formulation: Power

- Putting this together for Gaussian noise:

$$\Lambda \propto \int_{\Sigma} \exp(h_{\Sigma} \cdot \mathbf{s} - \mathbf{s} \cdot \mathbf{s}/2) d\mathbf{s}$$



# Formulation: Power

- Putting this together for Gaussian noise:

$$\Lambda \propto \int_{\Sigma} \exp(h_{\Sigma} \cdot \mathbf{s} - \mathbf{s} \cdot \mathbf{s}/2) d\mathbf{s}$$

- Note that  $\Lambda$  is monotonic in length of projection,  $|h_{\Sigma}|$ , so  $|h_{\Sigma}|$  is equivalent statistic.
- **Power** of projected data,  $|h_{\Sigma}|^2$ , is also equivalent optimal statistic and is most commonly used.



# Formulation: Power<sup>2.0</sup>

- For more than one detector, simply substitute:

$$h = A^k(\Omega) h_k(t)$$

Detector  
Antenna  
Pattern

Detector  
Data

# Formulation: Power<sup>2.0</sup>

- For more than one detector, simply substitute:

$$h = A^k(\Omega) h_k(t)$$

- Eg, for two detectors with Gaussian noise, optimal statistic is

$$\begin{aligned} |h_\Sigma|^2 &= | [A_1 h_1 + A_2 h_2]_\Sigma |^2 \\ &= A_1^2 |h_1^2|_\Sigma + A_2^2 |h_2^2|_\Sigma \\ &\quad + 2A_1 A_2 [h_1 \cdot h_2]_\Sigma \end{aligned}$$

# Formulation: Power<sup>2.0</sup>

- For more than one detector, simply substitute:

$$h = A^k(\Omega) h_k(t)$$

- Eg, for two detectors with Gaussian noise, optimal statistic is

$$\begin{aligned} |h_\Sigma|^2 &= | [A_1 h_1 + A_2 h_2]_\Sigma |^2 \\ &= \boxed{A_1^2 |h_1^2|_\Sigma + A_2^2 |h_2^2|_\Sigma} \quad \text{auto power} \\ &\quad + 2A_1 A_2 [h_1 \cdot h_2]_\Sigma \end{aligned}$$

# Formulation: Power<sup>2.0</sup>

- For more than one detector, simply substitute:

$$h = A^k(\Omega) h_k(t)$$

- Eg, for two detectors with Gaussian noise, optimal statistic is

$$\begin{aligned} |h_\Sigma|^2 &= | [A_1 h_1 + A_2 h_2]_\Sigma |^2 \\ &= \boxed{A_1^2 |h_1^2|_\Sigma + A_2^2 |h_2^2|_\Sigma} \quad \text{auto power} \\ &\quad + \boxed{2A_1 A_2 [h_1 \cdot h_2]_\Sigma} \quad \text{cross power} \end{aligned}$$

# Application: we know $s$ .

- For the cases where we know the signal  $\hat{s}$  we are looking for, the measure becomes a Dirac delta:

$$\mathcal{D}[s] = \delta(s - \hat{s}) ds$$

- Integrating against this measure,  $h_{\Sigma} = h \cdot \hat{s} / |\hat{s}|^2$ .
- This is nothing but the *signal-to-noise* ratio of the matched filter search.
- If you are searching for one of a finite set of signals, search for each individually - template bank.
- This is how we search for binaries and pulsars.

# Application: $s$ is random.

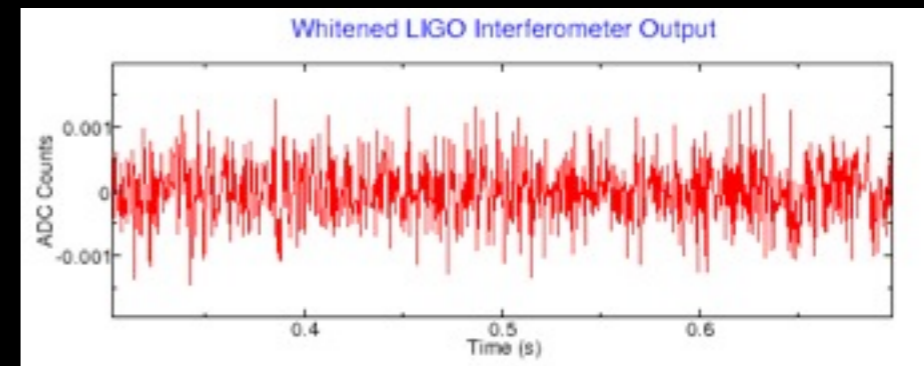
- For stochastic backgrounds, by central limit theorem, resulting signal is random and Gaussian.
- Sum of Gaussians,  $h = n + s$ , is Gaussian.  $h_{\Sigma} = h$ .
- If  $\sigma_n^2 \gg \sigma_s^2$  then  $h_{\Sigma} \cdot h_{\Sigma} \approx n \cdot n$  regardless of whether signal exists, so can't use auto-power.
- However, for two instruments with data  $h_1$  and  $h_2$  of length  $N$ , noises  $n_1$  and  $n_2$  are uncorrelated, so  $n_1 \cdot n_2 \sim \sqrt{N}$ , but  $s_1$  and  $s_2$  are correlated, so  $s_1 \cdot s_2 \sim N$ .
- So, in the limit of large enough  $N$ ,  $h_1 \cdot h_2 \sim s_1 \cdot s_2$ .

# Application: $s$ not known.

- For a determinate signal that is not completely known we can again apply likelihood.
- In this case,  $\mathcal{D}[s]$  encodes what we know about the signal.
- Eg, to search for signals that last  $\Delta t$  seconds, use data segments  $h_{\Sigma}$  of that duration. If additionally signal frequencies are known to be in band  $\Delta f$ , use data segments  $h_{\Sigma}$  restricted to that band.
- So, for a single detector, an optimal statistic is auto-power for data segments restricted to  $\Delta t \Delta f$ .

# Application: $s$ not known.

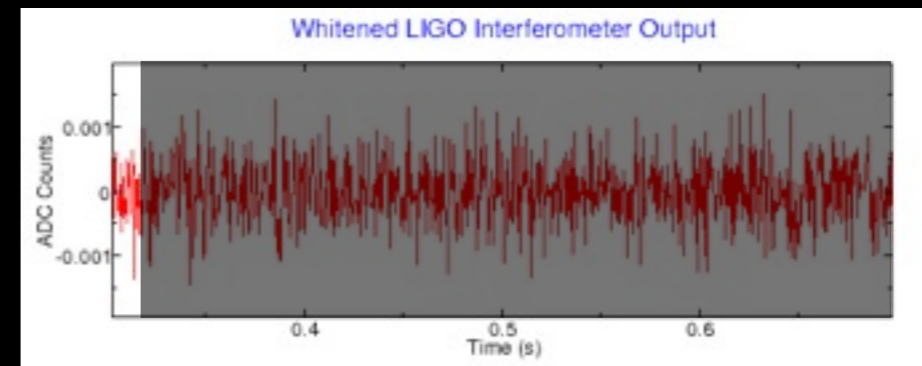
- How do we practically restrict to data segments of dimension  $\Delta t \Delta f$ ?
  - take a slice of detector data and Fourier transform it.





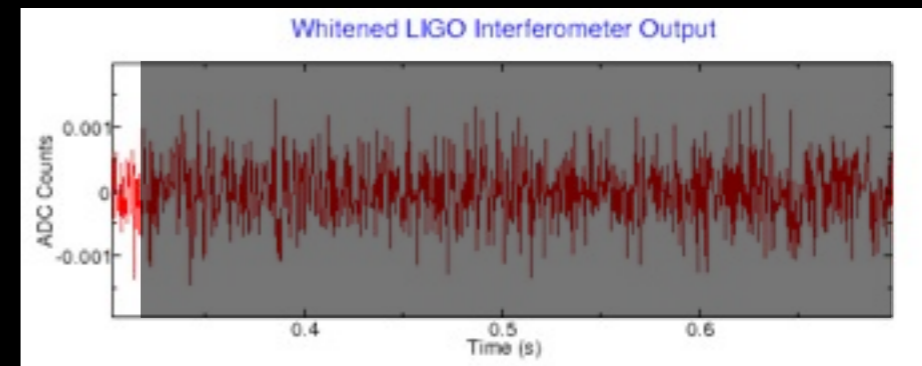
# Application: $s$ not known.

- How do we practically restrict to data segments of dimension  $\Delta t \Delta f$ ?
  - take a slice of detector data and Fourier transform it.



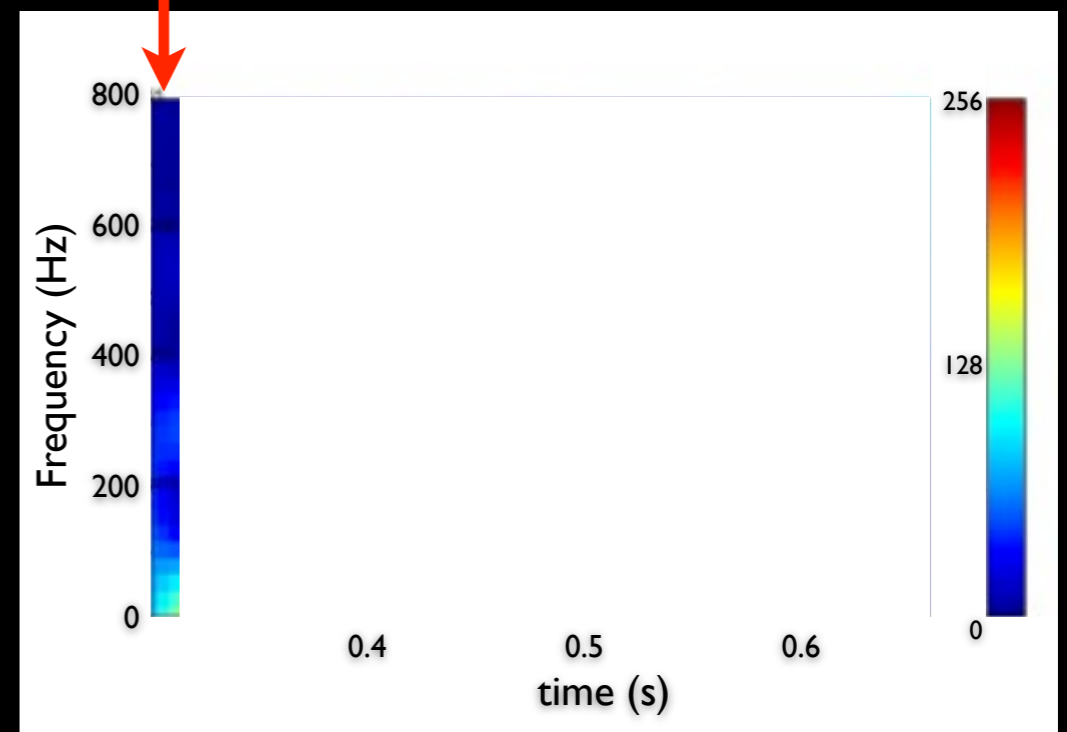
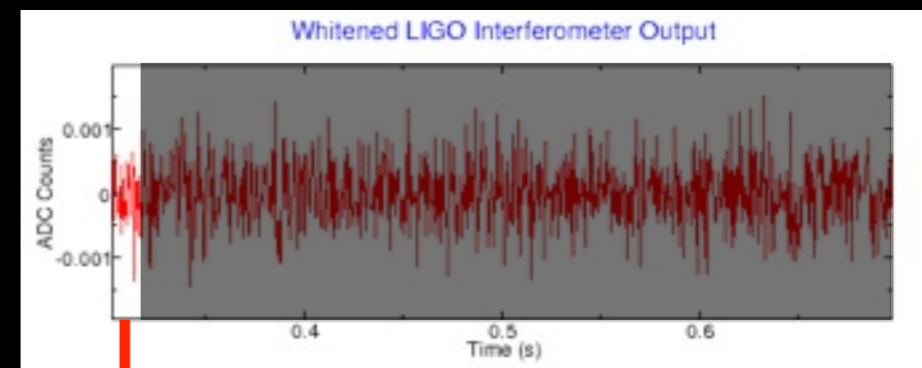
# Application: $s$ not known.

- How do we practically restrict to data segments of dimension  $\Delta t \Delta f$ ?
  - take a slice of detector data and Fourier transform it.
  - plot the Fourier coefficient magnitudes on a vertical line.



# Application: $s$ not known.

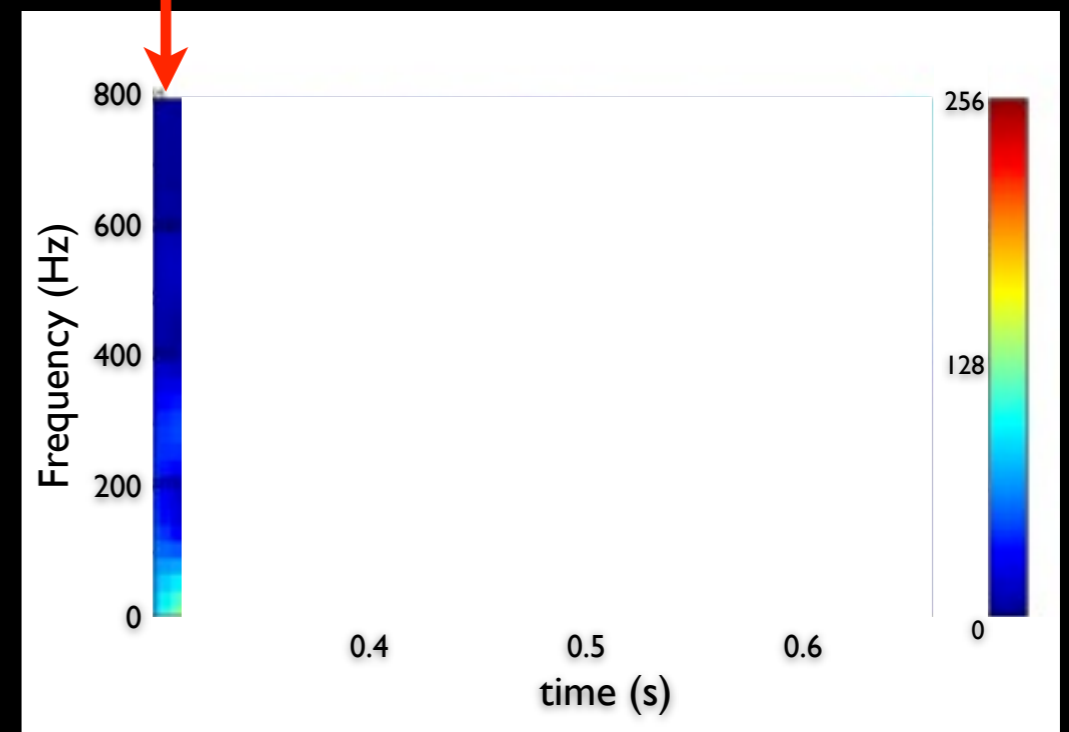
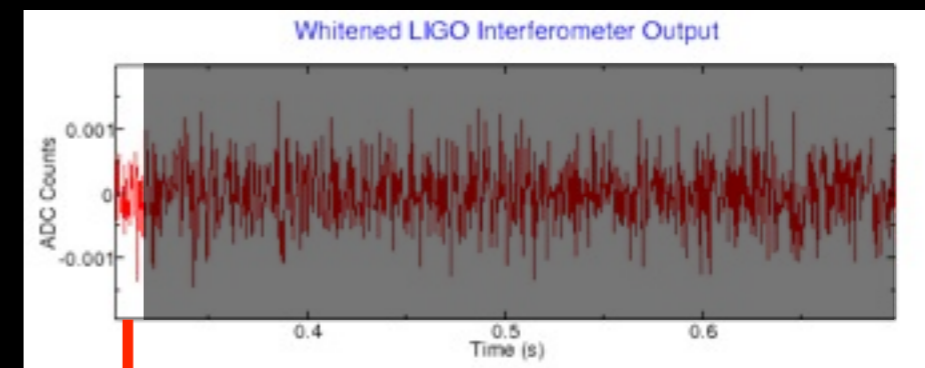
- How do we practically restrict to data segments of dimension  $\Delta t \Delta f$ ?
  - take a slice of detector data and Fourier transform it.
  - plot the Fourier coefficient magnitudes on a vertical line.



# Application: $s$ not known.

- How do we practically restrict to data segments of dimension  $\Delta t \Delta f$ ?

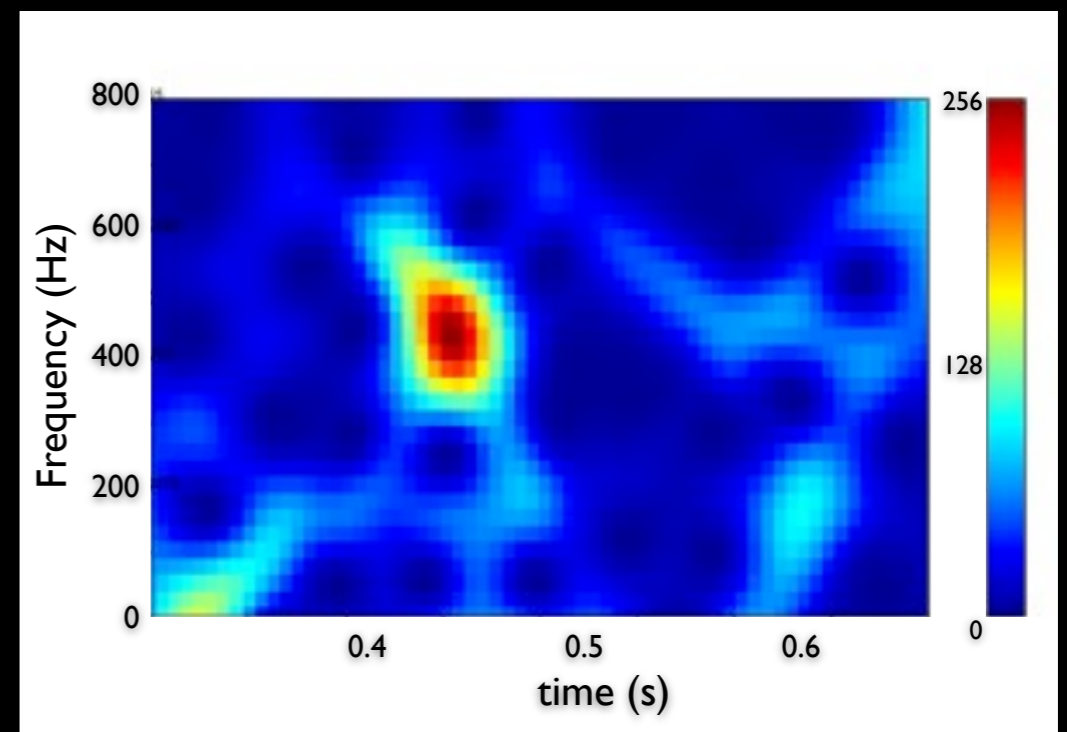
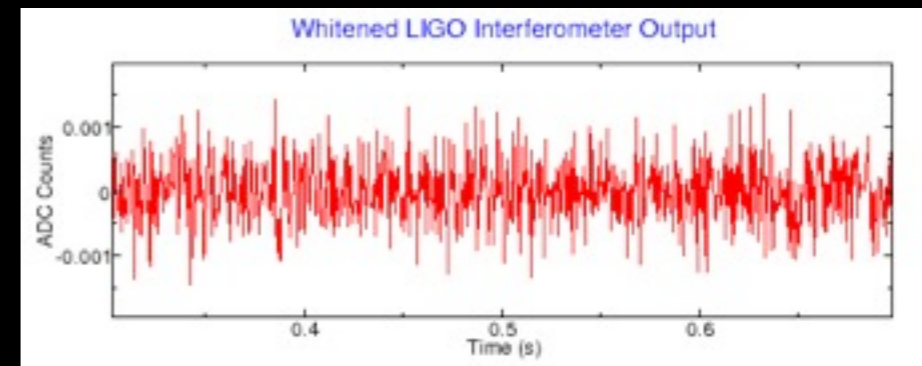
- take a slice of detector data and Fourier transform it.
- plot the Fourier coefficient magnitudes on a vertical line.
- repeat for subsequent slices of data.



# Application: $s$ not known.

- How do we practically restrict to data segments of dimension  $\Delta t \Delta f$ ?

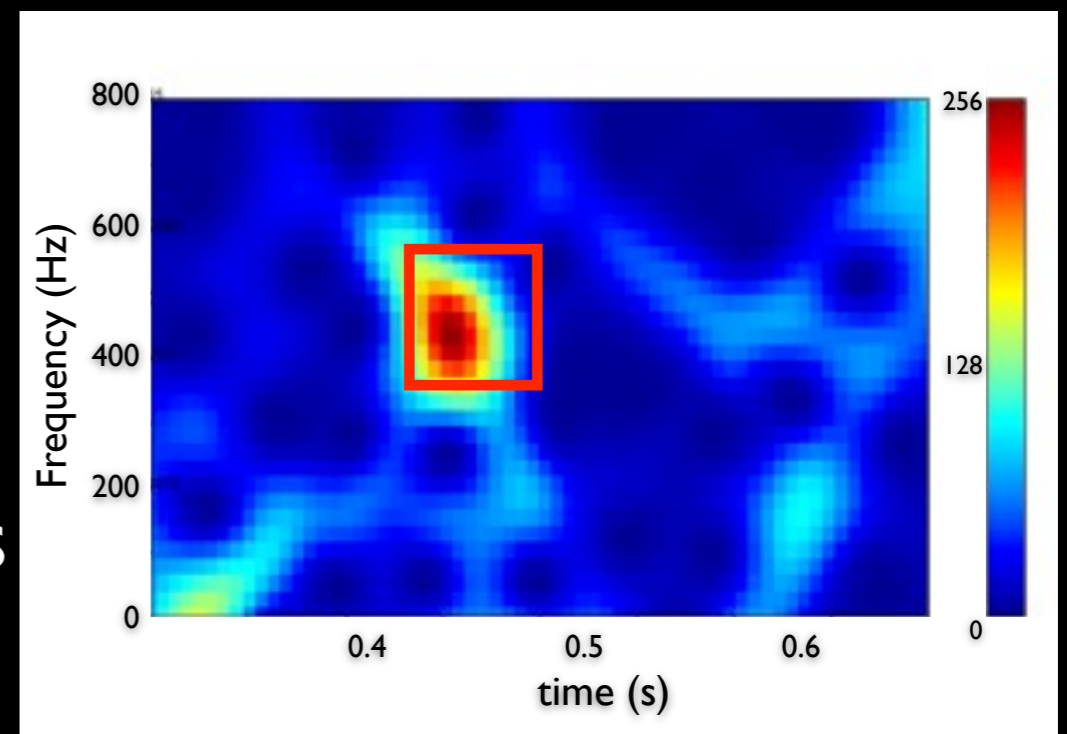
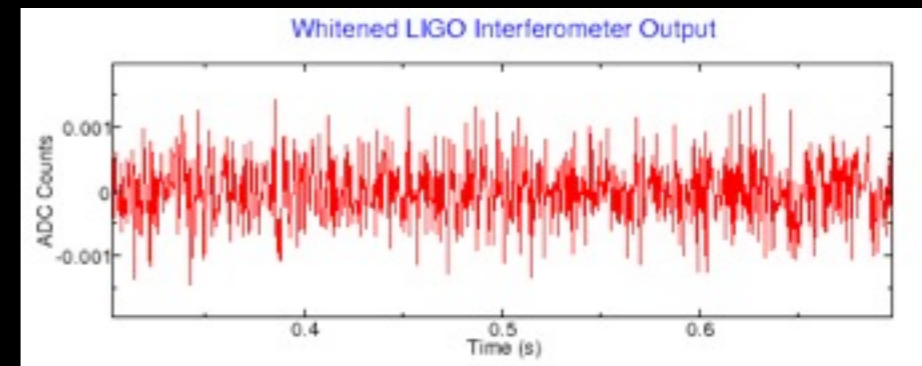
- take a slice of detector data and Fourier transform it.
- plot the Fourier coefficient magnitudes on a vertical line.
- repeat for subsequent slices of data.



# Application: $S$ not known.

- How do we practically restrict to data segments of dimension  $\Delta t \Delta f$ ?

- take a slice of detector data and Fourier transform it.
- plot the Fourier coefficient magnitudes on a vertical line.
- repeat for subsequent slices of data.
- then we can search for boxes with statistical significance.



# Application: real life

- In an ideal world, these four data analysis groups and these methods would be the end of the story, but ...
- In reality:
  - signals are never exactly known
  - no theoretical model for noise statistics
  - false alarm probabilities can't be calculated
  - narrowband noise “lines” can mask signals
  - “glitches” (burst of noise) mimic real signals
  - noise not really uncorrelated between detectors
  - ...

# Migration: burst? glitch?

- Consider the problem looking for bursts in real noise.
- Power no longer optimal because loud glitches also cause large auto-power.
- Question: Who knows how to search for signals when you can't tell the signal from a detector's noise?
- Answer: Analysts who look for stochastic backgrounds!

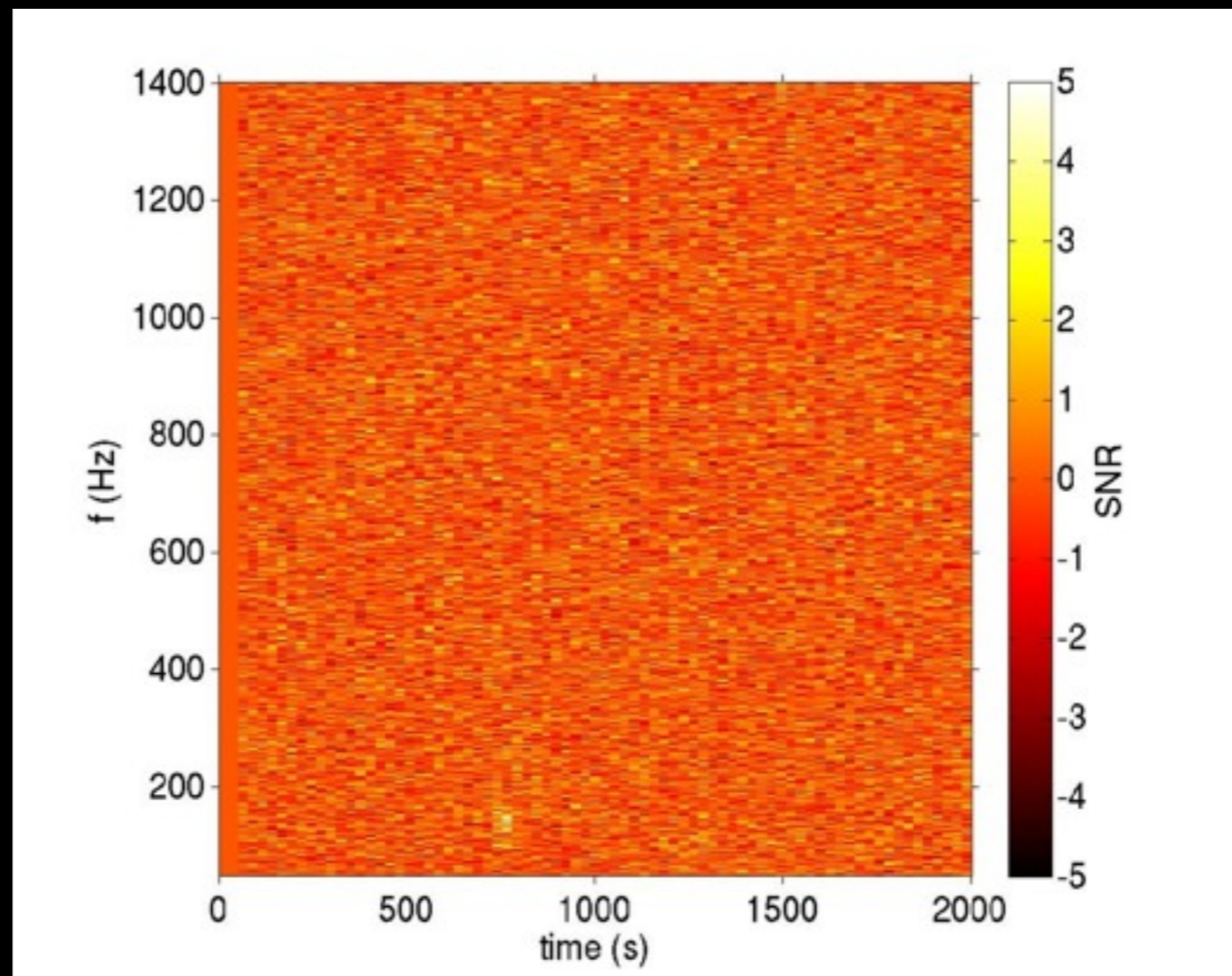


# Migration: STAMP

- Idea: Use cross-power to look for unmodeled burst signals in glitchy data.
- Led by: LIGO stochastic analysis group
- Called: STAMP - Stochastic Transient Analysis Multi-detector Pipeline.
- Uses: analysis code and expertise from
  - stochastic analyses
  - unmodeled burst analyses
  - pulsar analyses

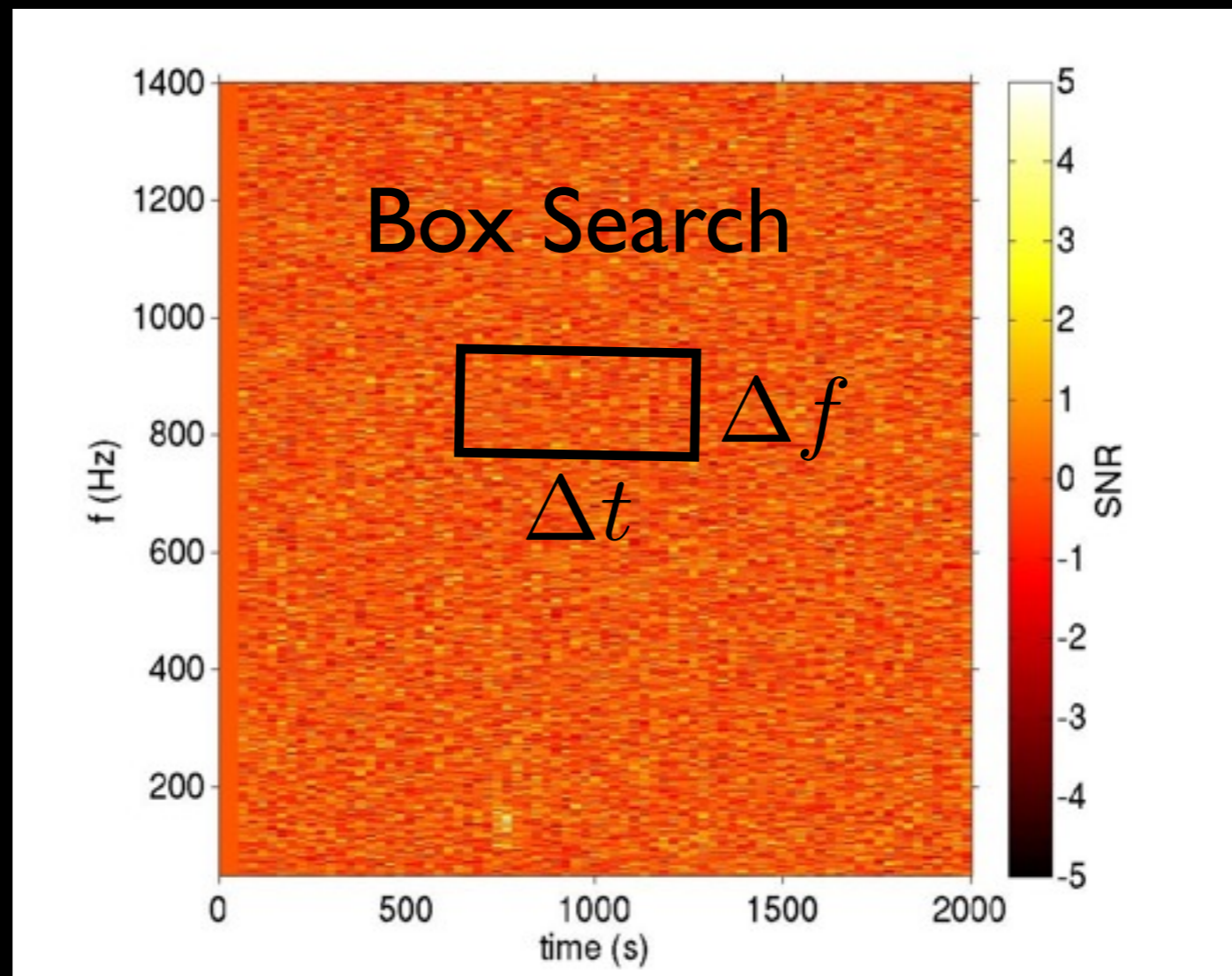
# Migration: how it works

- Uses a TF representation of cross-power to project onto signal space.



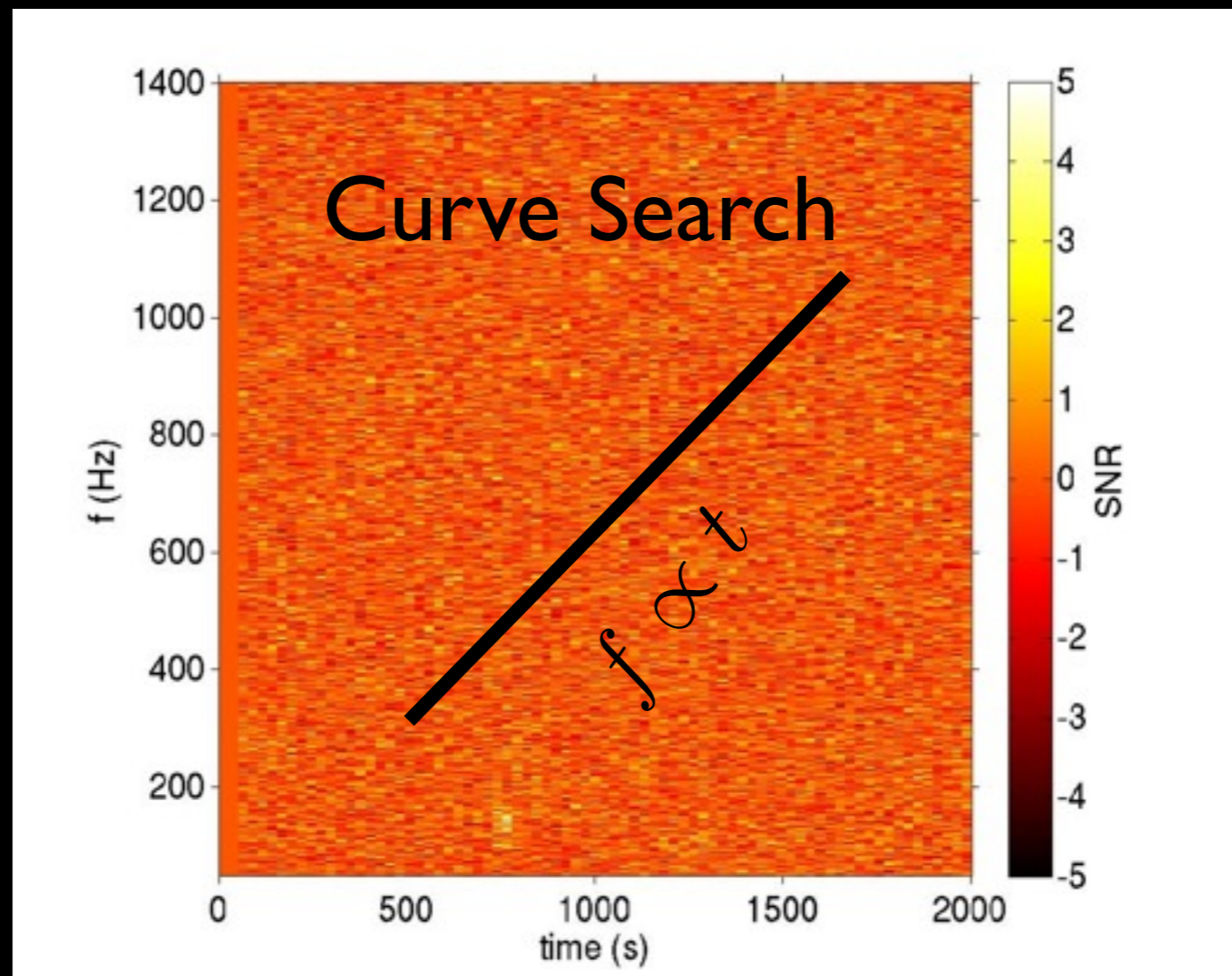
# Migration: how it works

- Uses a TF representation of cross-power to project onto signal space.



# Migration: how it works

- Uses a TF representation of cross-power to project onto signal space.



# Migration: activity

- So far, the STAMP group have produced:
  - a methods paper (*PRD 83, 083004*)
  - a detector noise paper (*CQG 28, 235008*)
  - a long GRB upper limits paper (*PRD 88, 122004*)
- Group is currently working on:
  - an all-sky search
  - an neutron star r-mode search
- Contributors include:
  - current** - Marie Anne Bizouard, Samuel Franco, Patrice Hello, Nelson Christensen, Eric Thrane, Shivaraj Kandhasamy, Tanner Prestegard, Patrick Meyers, Jialun Luo, Michael Coughlin, Bernard Whiting, Antonis Mytidis.
  - past** - Christian Ott, Steven Dorsher, Stefanos Giampanis, Vuk Mandic, Peter Raffai, WGA

# Conclusion

- All “four types” of LIGO-Virgo data analysis have a lot in common:
  - use the same data
  - based on likelihood
  - the only difference is signal space we project on.
- Expertise and tools are portable across many analyses.
- STAMP is proving to be an example of fruitful interactions between analysts from different camps.