

# The Time Domain



# The Sample Interval

We have a sequence of  $N$  observations

$$z_n, \quad n = 0, 1, 2, \dots N - 1$$

which coincide with times

$$t_n, \quad n = 0, 1, 2, \dots N - 1.$$

The sequence  $z_n$  is called a *discrete time series*.

It is assumed that the *sample interval*,  $\Delta$ , is constant

$$t_n = n\Delta$$

with the time at  $n = 0$  defined to be 0. The *duration* is  $T = (N - 1)\Delta$ .

If the sample interval in your data is not uniform, the first processing step is to interpolate it be so.



# The Underlying Process

A critical assumption is that there exists some “process”  $z(t)$  that our data sequence  $z_n$  is a *sample of*:

$$z_n = z(n\Delta), \quad n = 0, 1, 2, \dots N - 1.$$

Unlike  $z_n$ ,  $z(t)$  is believed to exist for *all times*.

- (i) The process  $z(t)$  exists in *continuous time*, while  $z_n$  only exists at *discrete times*.
- (ii) The process  $z(t)$  exists for *all past and future times*, while  $z_n$  is only available over a certain time interval.

It is the properties of  $z(t)$  that we are trying to estimate, *based on* the available sample  $z_n$ .



# Measurement Noise

In reality, the measurement device and/or data processing probably introduces some artificial variability, termed *noise*.

It is more realistic to consider that the observations contain samples of the process of interest,  $z(t)$ , *plus* some noise  $\epsilon_n$ :

$$z_n = z(n\Delta) + \epsilon_n, \quad n = 0, 1, 2, \dots N - 1.$$

This is an example of the *unobserved components model*. This means that we *believe* that the data is composed of *different components*, but we cannot observe these components individually.

The process  $z(t)$  is potentially obscured or degraded by the limitations of data collection in three ways: (i) finite sample interval, (ii) finite duration, (iii) noise.

Because of this, the time series is an *imperfect* representation of the real-world processes we are trying to study.



# A Pair of Time Series

In oceanography we often have a *pair* of time series  $x_n$  and  $y_n$ . Such data is called *bivariate*, meaning that it consists of two variables.

These may represent horizontal velocity (as in current meters) or displacement (floats or drifters).

Bivariate data can be thought of as a vector having two elements:

$$\mathbf{z}_n = \begin{bmatrix} x_n \\ y_n \end{bmatrix}.$$

The subscript  $n$  here refers to  $n$  different copies of the vector, *not* to the elements of that vector!

Alternatively, we can also think of this data consisting of a single *complex-valued* time series  $z_n \equiv x_n + iy_n$ , where  $i \equiv \sqrt{-1}$ .

Vector and complex notations will both be discussed in detail later.



# Time versus Frequency

There are two opposite points of view regarding the time series  $z_n$ .

The first regards  $z_n$  as being built up as a sequence of discrete values  $z_0, z_2, \dots, z_{N-1}$ .

This is the domain of *statistics*: the mean, variance, histogram, etc.

When we look at data statistics, generally, the order in which the values are observed *doesn't matter*.

The second point of view regards  $z_n$  as being built up of sinusoids: purely periodic functions spanning the whole duration of the data.

This is the domain of *Fourier spectral analysis*.

In between these two extremes is wavelet analysis.

This lecture will focus on what can be done in the time domain.



# Time-Domain Statistics

A good place to start is with the very simplest tools. We'll change to  $x_n$  and  $x(t)$  as this discussion pertains to real-valued data.

The *sample mean* describes a “typical” value:

$$\bar{x} \equiv \frac{1}{N} \sum_{n=0}^{N-1} x_n$$

The *sample variance* gives the spread about the mean:

$$\sigma_x^2 \equiv \frac{1}{N} \sum_{n=0}^{N-1} (x_n - \bar{x})^2$$

“Sample” here means that it is computed from the observed data, as opposed to being a property of the assumed underlying process  $x(t)$ .

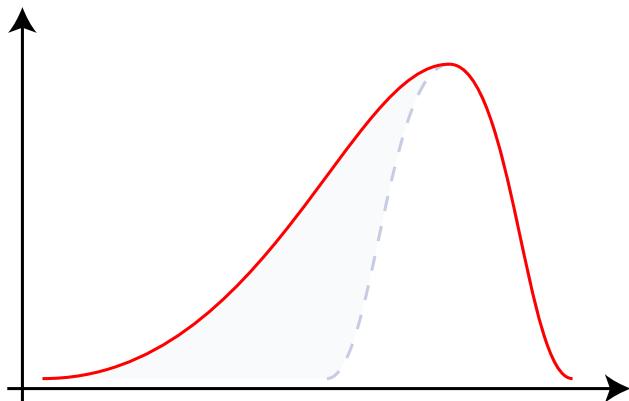


- The mean and variance are called the first two *moments* of the distribution of values associated with the process.

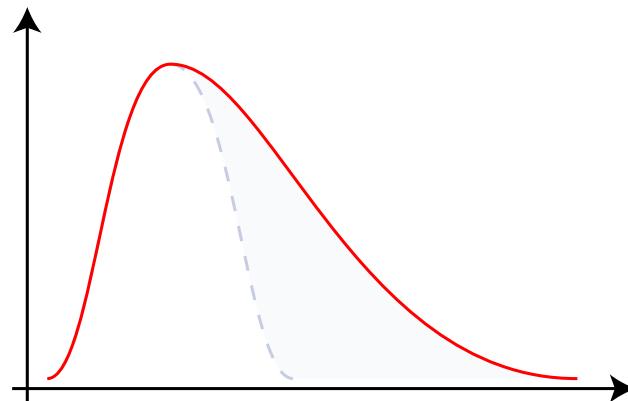
# Skewness

The *skewness* describes the tendency for an *asymmetry* between positive excursions and negative excursions:

$$\gamma_x \equiv \frac{1}{\sigma_x^3} \frac{1}{N} \sum_{n=0}^{N-1} (x_n - \bar{x})^3$$



Negative Skew



Positive Skew

# Kurtosis

The *kurtosis* is said to either measure *peakedness* (concentration near  $\bar{x}$ ), or a tendency for *long tails* (concentration far from  $\bar{x}$ ):

$$\kappa_x \equiv \frac{1}{\sigma_x^4} \frac{1}{N} \sum_{n=0}^{N-1} (x_n - \bar{x})^4$$

Actually, it measures both. Kurtosis is a measure of the spread of  $x_n$  about the *two points*  $\bar{x} \pm \sigma_x$ . This can happen *either* for peakness or for long tails! *See Moors (1986), “The Meaning of Kurtosis”.*

Because the value of kurtosis for a Gaussian process can be shown to be equal to 3, one sometimes encounters the *excess kurtosis*

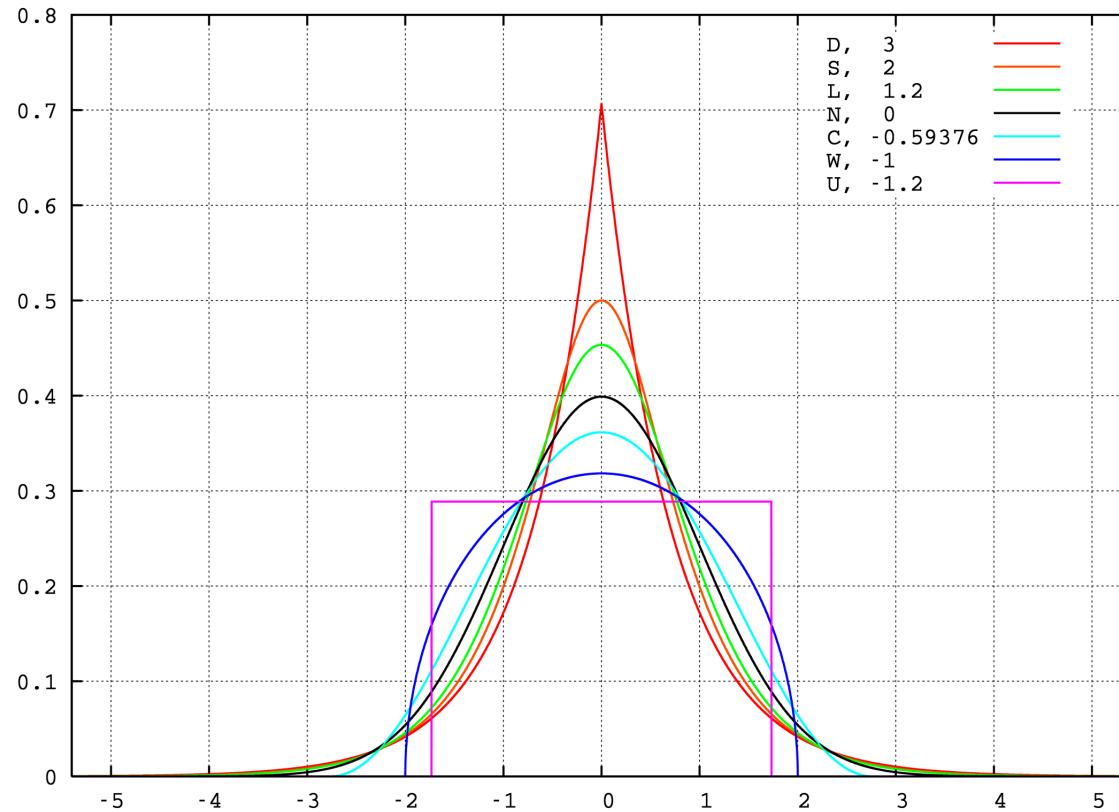
$$\tilde{\kappa}_x \equiv \kappa_x - 3.$$

Values of  $\tilde{\kappa}_x > 0$  mean the data is *more kurtotic*—peaked or long-tailed—than a Gaussian, while  $\tilde{\kappa}_x < 0$  means it is less so.



# Illustration of Kurtosis

Distributions corresponding to different values of excess kurtosis.

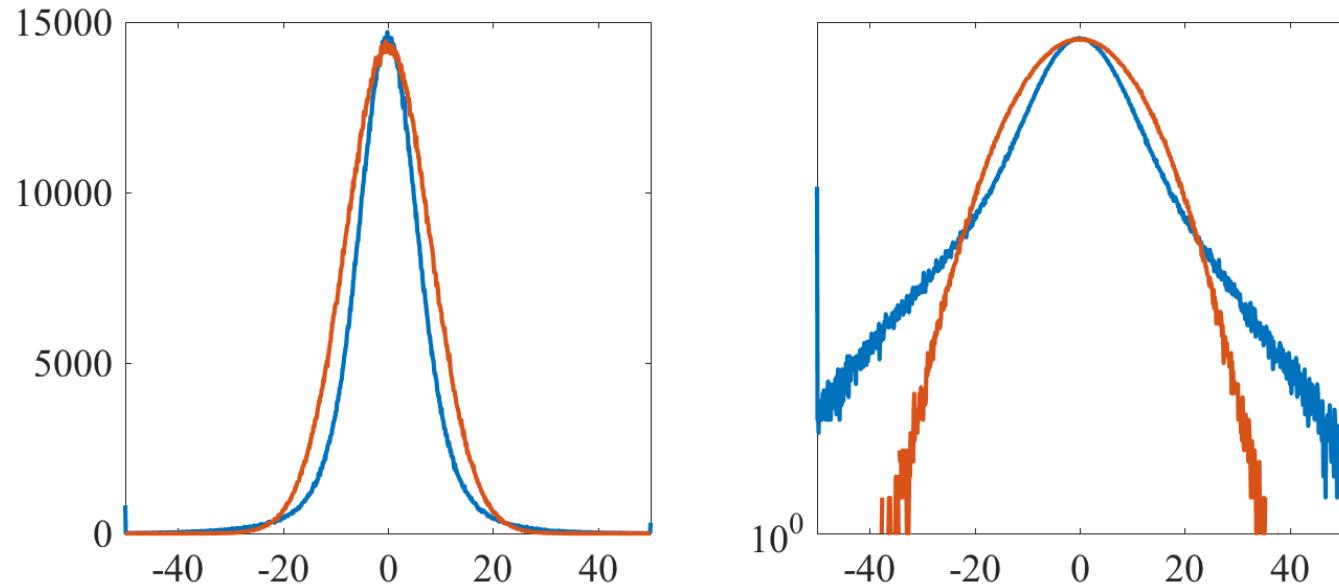


Positive excess kurtosis corresponds to long tails and peakedness.



# Histogram

The mean, variance, skewness, and kurtosis describe aspects of the *histogram*: the observed distribution of data values.

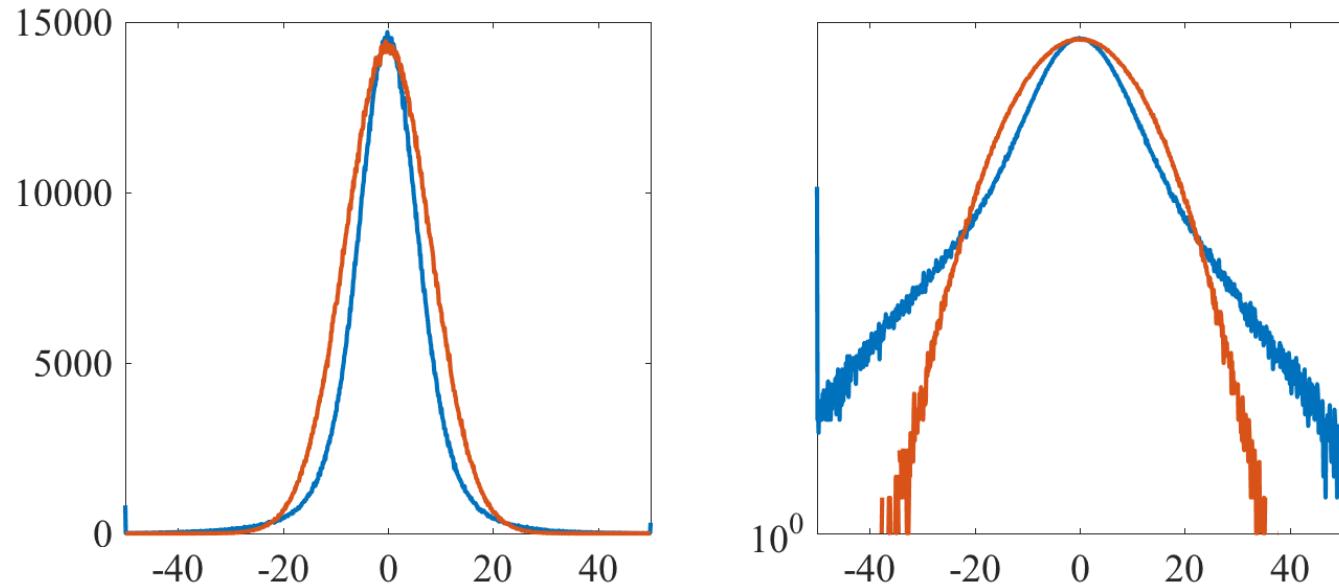


Here is the histogram of *all* SSH values from long altimeter track (blue), versus Gaussian noise having the same variance (orange).



# Histogram

The mean, variance, skewness, and kurtosis describe aspects of the *histogram*: the observed distribution of data values.



Here is the histogram of *all* SSH values from long altimeter track (blue), versus Gaussian noise having the same variance (orange).



# Simple Smoothing

One of the most effective ways to process a time series is with a simple smoothing.

Let  $g_m$  be a length  $M$  sequence, where  $M$  is *odd*, defined for

$$-(M-1)/2, \dots, -2, -1, 0, 1, 2, \dots, (M-1)/2.$$

Note that we define  $g_m$  to be centered on  $m = 0$ , instead of running between 0 and  $M - 1$ .

A *smoothed* version of the discrete time series  $z_n$  is defined as

$$\tilde{z}_n = \sum_{m=-(M-1)/2}^{(M-1)/2} z_{n-m} g_m$$

where  $g_m$  is called the *filter* or the *smoothing window*. It is also useful to examine the *residuals* from the original,  $\check{z}_n \equiv z_n - \tilde{z}_n$ .



# Simple Smoothing Example

An example of simple smoothing is a *running mean*. A five-point running mean is given by:

$$\tilde{z}_n = \frac{1}{5} [z_{n-2} + z_{n-1} + z_n + z_{n+1} + z_{n+2}].$$

This is expressed by the filtration equation

$$\tilde{z}_n = \sum_{m=-(M-1)/2}^{(M-1)/2} z_{n-m} g_m$$

with the choice

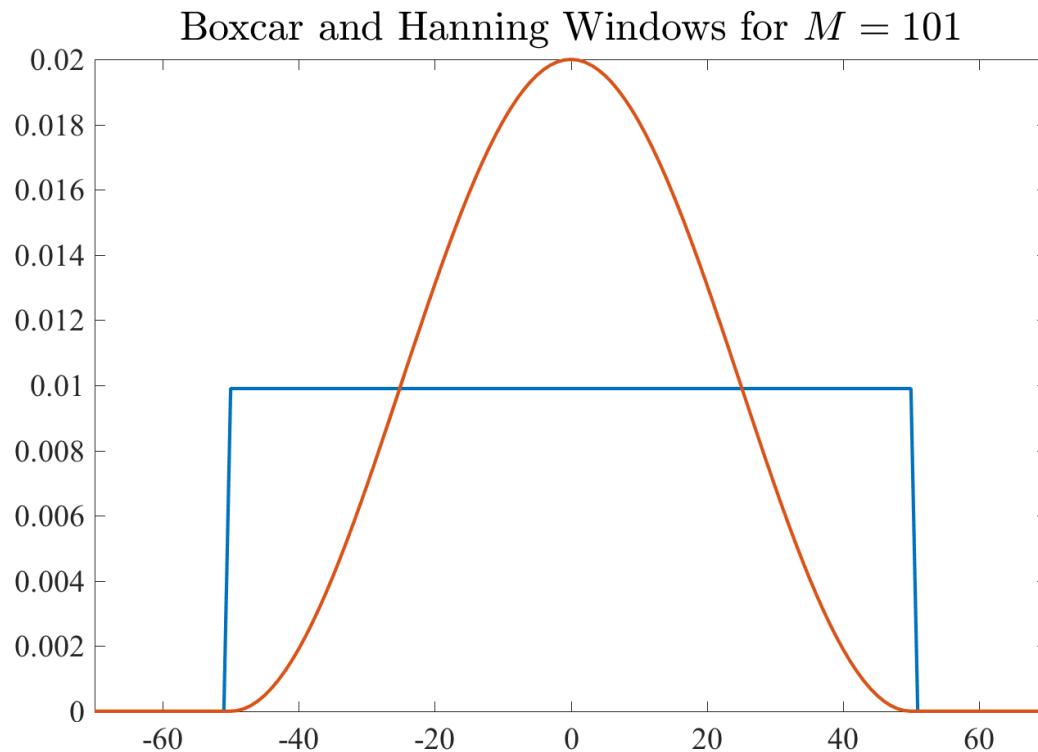
$$g_m = 1/5, \quad m = -2, -1, 0, 1, 2.$$

The simplest choice of filter is  $g_m = 1/M$ , a constant over the  $M$  points. Then the filtration defines an  $M$ -point running mean.



# Choice of Filter

The running mean filter  $g_m = 1/M$  is called the *boxcar* or *rectangle function*. Another popular choice is the *Hanning window*.



The Hanning window is just a half-period of a cosine, offset.

# How to Choose a Filter

The goal of simple smoothing is to separate relatively “fast” from relatively “slow” variability.

Many functions can be used as smoothing filters. However, for a first look at the data, the details of the filter are not so important.

The important thing is to define a sensible *weighted average*.

The boxcar filter has sharp “edges” that can lead to artifacts, as we will see later. Also, the boxcar is highly distributed, and doesn’t place emphasis on the “present time” compared to nearby times.

For these reasons, the Hanning window is sometimes more appropriate for simple smoothing.

In jLab simple smoothing is carried out with `vfilt`.



# What to do at Endpoints?

Smoothing runs into a difficulty near the endpoints of  $z_n$ :

$$\tilde{z}_n = \sum_{m=-(M-1)/2}^{(M-1)/2} z_{n-m} g_m.$$

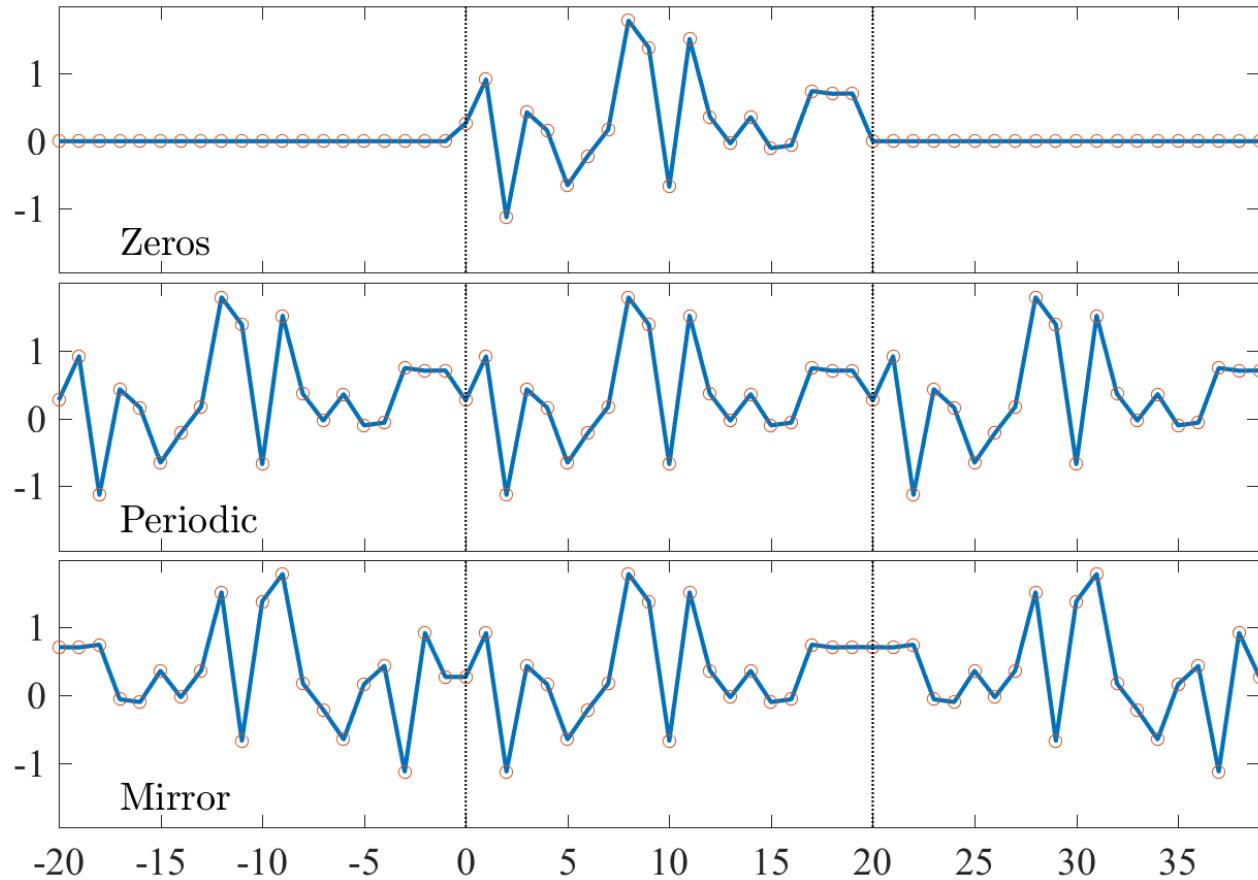
When we are within a filter half-width  $(M - 1)/2$  of the beginning or end of  $z_n$ , the filter “falls off” the end of the data.

Some choice must be made in order to have the smoothed version  $\tilde{z}_n$  of the data be well defined. There are five common choices.

1. **Truncate:** Omit affected points, such that the length of  $\tilde{z}_n$  will be about  $M$  points *less than* the length of  $z_n$ .
2. **NANs:** Replace these with NaNs or *indeterminate* values.
3. **Zeros:** Set  $z_n$  equal to zero for  $n \leq 0$  or  $n \geq N - 1$ .
4. **Periodic:** Make  $z_n$  periodic by wrapping around the ends.
5. **Mirror:** Reflect  $z_n$  about its beginning and also about its end.



# Endpoint Illustration



The *mirror* condition generally leads to the fewest “edge effects”, especially when the data is nonstationary or has a linear trend.



# Summary

This lecture has focused on



# Summary

This lecture has focused on

- Introducing the concepts of *discrete sampling*, *sample interval*, *measurement noise*, and the *underlying process*.



# Summary

This lecture has focused on

- Introducing the concepts of *discrete sampling*, *sample interval*, *measurement noise*, and the *underlying process*.
- Defining the first four *moments*—mean, variance, skewness, and kurtosis—as well as the *histogram*.



# Summary

This lecture has focused on

- Introducing the concepts of *discrete sampling*, *sample interval*, *measurement noise*, and the *underlying process*.
- Defining the first four *moments*—mean, variance, skewness, and kurtosis—as well as the *histogram*.
- Discussing *simple smoothing* and details of its implementation.



# Summary

This lecture has focused on

- Introducing the concepts of *discrete sampling*, *sample interval*, *measurement noise*, and the *underlying process*.
- Defining the first four *moments*—mean, variance, skewness, and kurtosis—as well as the *histogram*.
- Discussing *simple smoothing* and details of its implementation.

My experience is that *looking at data* together with *statistics* and *simple smoothing* is maybe 50% of analyzing time series!

There are more sophisticated tools that can often, but not always, be very useful in unlocking the potential of the data.

However, learning how to make use of these takes a lot more work!

To be continued...



# Three Cases

In general we have three types of data.

1. Gridded data
2. Dense irregular data
3. Sparse irregular data

We have to approach these differently.



# Gridded data

With gridded data we can fruitfully analyze using two-dimensional statistics by directly averaging in different directions along a 2D, 3D, or N-D “cube” of data.

Averaging can be done quickly, without explicit loops, for such data. Then we imagine turning the cube in different directions and averaging along different axes.

It is often useful to split time into two dimensions. For example,

$\text{lat} \times \text{lon} \times \text{time}$

can be reorganized to become

$\text{lat} \times \text{lon} \times \text{time of year} \times \text{different years.}$

So then averaging over the 4th dimension creates a composite year, while taking the standard deviation gives the year-to-year variability.



For this, `reshape` and `permute` are useful.

# Gridded data

If the data is so big that we can't load it into memory all at one, then we can average by aggregating: loading one time slice in at a time, adding to a running total, then dividing by the number of slices at the end.

The variance can also be computed in this way using formulas like

$$\overline{(x - \bar{x})^2} = \overline{x^2} - \bar{x}^2.$$

So we can create the variance by creating aggregated averages of both  $\bar{x}$  and  $\overline{x^2}$ . The same applies to higher-order moments.

We can also use `twodhist` and `twodstats` to examine distributions and averages in parameter space, e.g. the distribution of sea surface temperature vs. sea level pressure for all locations and all times.

If we find a pattern in parameter space, it can be quantified using regression analysis, if desired.



# Dense irregular data

For irregularly sampled (non-gridded) data, we cannot directly average, but we can use `twodhist` and `twodstats` to examine distributions and averages is sensible.

If there are enough data points such these histograms are sufficiently “filled in”, we will call data “dense”.

With dense irregular data, as with gridded data, we imagine the dataset to be a large multivariate “cloud”, e.g.

$\text{lat} \times \text{lon} \times \text{time} \times \text{temperature} \times \text{pressure}$

and then we use the distributional analysis to slice it in different ways, looking for patterns.

It is often very useful to have two-dimensional statistics with quantities of different units on the x- and y-axes, e.g. latitude  $\times$  time. Then we can look at distributions, means, and standard deviations, etc. in this plane.



# Sparse data

Sparse simply means there is not very much data. Practically speaking, our approach must be different when there is not enough data to fill in a histogram.

In this case, we learn a lot by employing scatter plots making creative use of size and/or color of symbols.

In Matlab this is done using scatter.



# Statistics Example

As an example of how to use time-domain statistics, we will look at a numerical model of the Gulf of Mexico.

First consider the mean and standard deviation of the velocity  $\mathbf{z}_n(x, y) \equiv [u_n(x, y) \ v_n(x, y)]^T$ ,

$$\bar{\mathbf{z}}(x, y) \equiv \frac{1}{N} \sum_{n=0}^{N-1} \mathbf{z}_n, \quad \sigma^2(x, y) \equiv \frac{1}{N} \sum_{n=0}^{N-1} (\mathbf{z}_n - \bar{\mathbf{z}})^T (\mathbf{z}_n - \bar{\mathbf{z}})$$

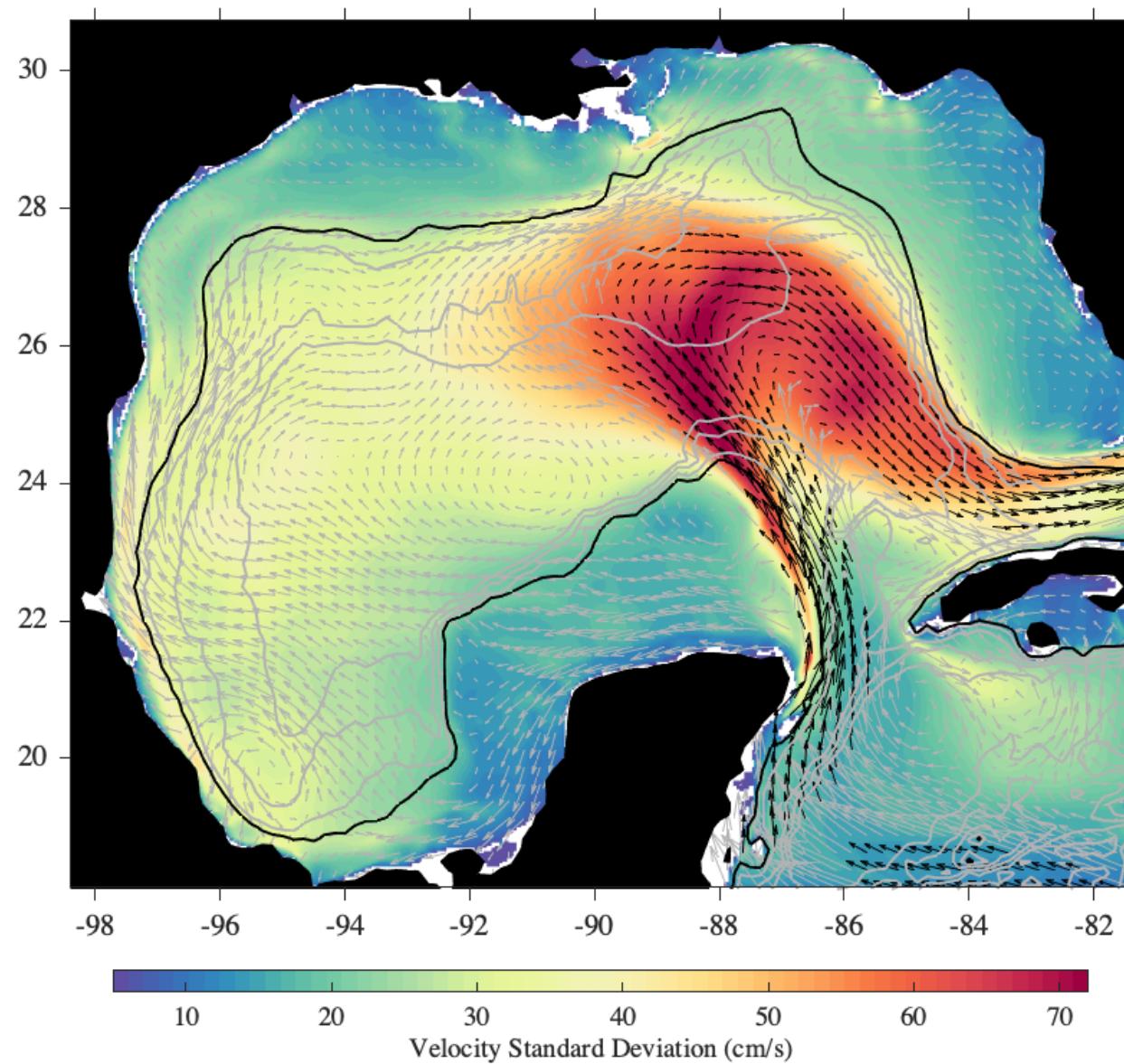
and the ratio of the mean flow magnitude to the standard deviation

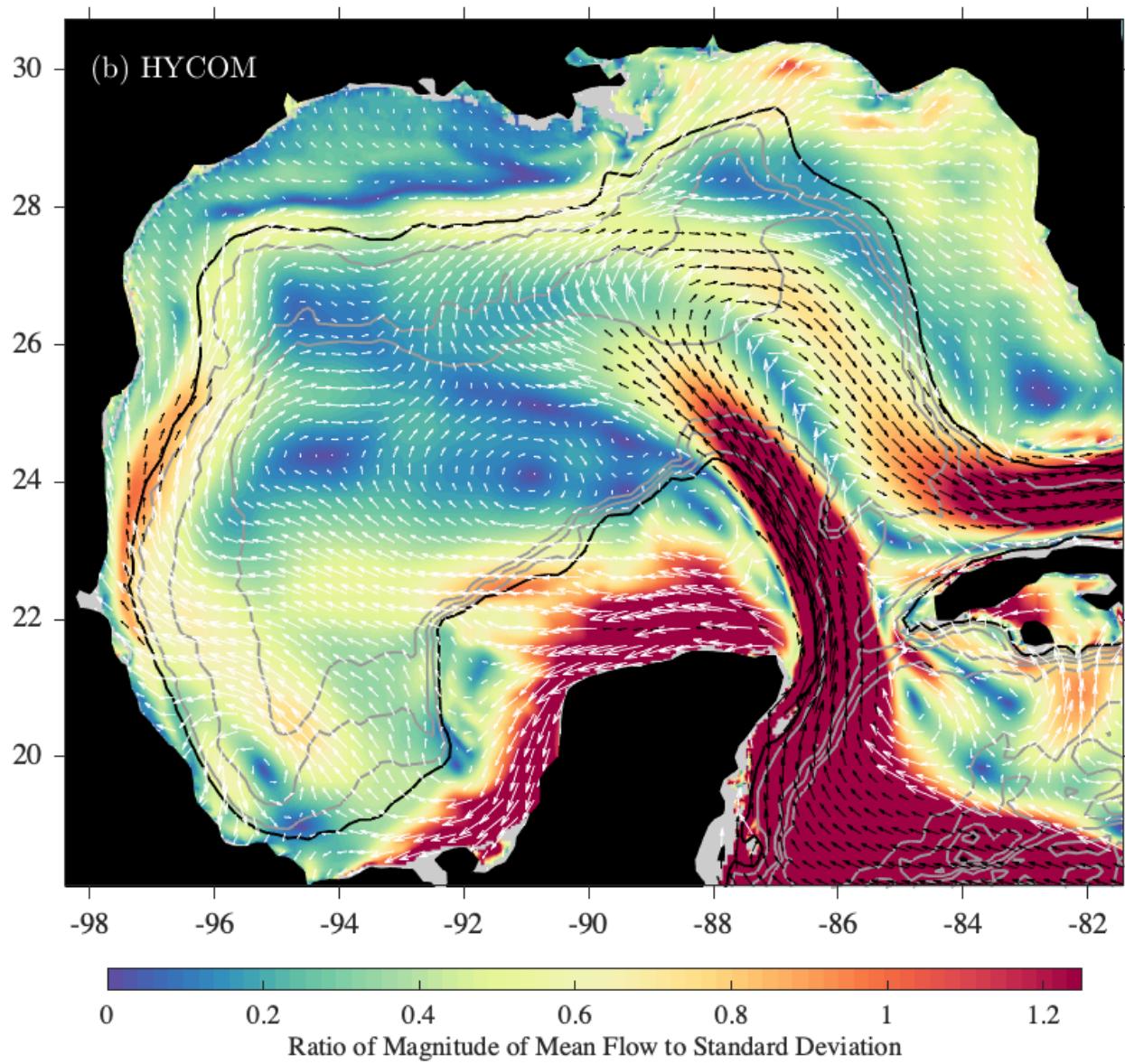
$$\frac{\|\bar{\mathbf{z}}(x, y)\|}{\sigma(x, y)} = \frac{\sqrt{\bar{u}^2(x, y) + \bar{v}^2(x, y)}}{\sigma(x, y)}$$

(with  $\|\mathbf{z}\| \equiv \sqrt{\mathbf{z}^T \mathbf{z}}$ ) which could be interpreted as a nondimensional measure of the *stability* of the flow patterns.



(Model courtesy of J. Zavala-Hidalgo and colleagues.)





# Statistics Example

Next we will look at the first three moments of the vorticity

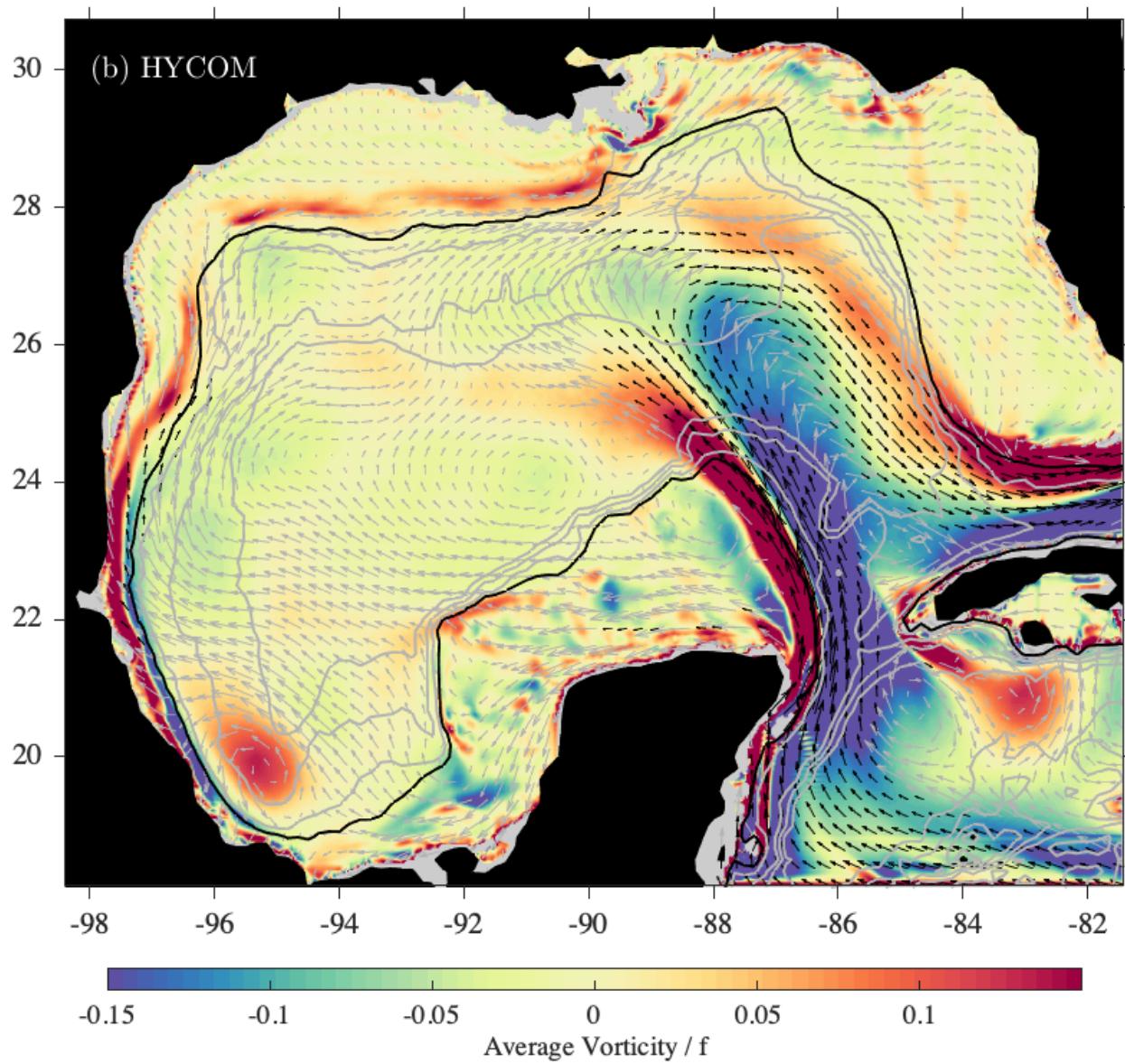
$$\zeta_n(x, y) \equiv \frac{\partial v_n}{\partial x} - \frac{\partial u_n}{\partial y}$$

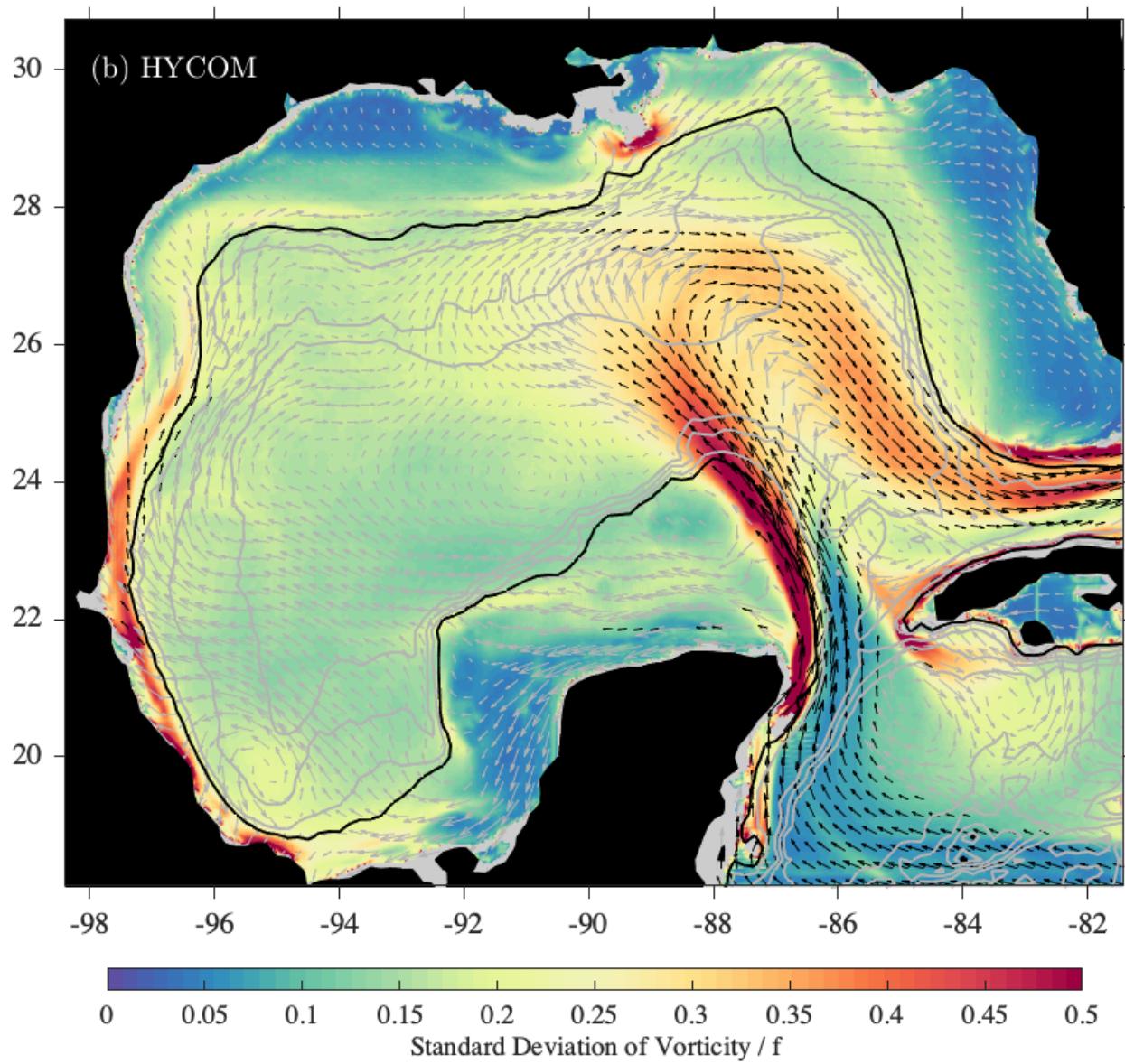
$$\bar{\zeta}(x, y) \equiv \frac{1}{N} \sum_{n=0}^{N-1} \zeta_n$$

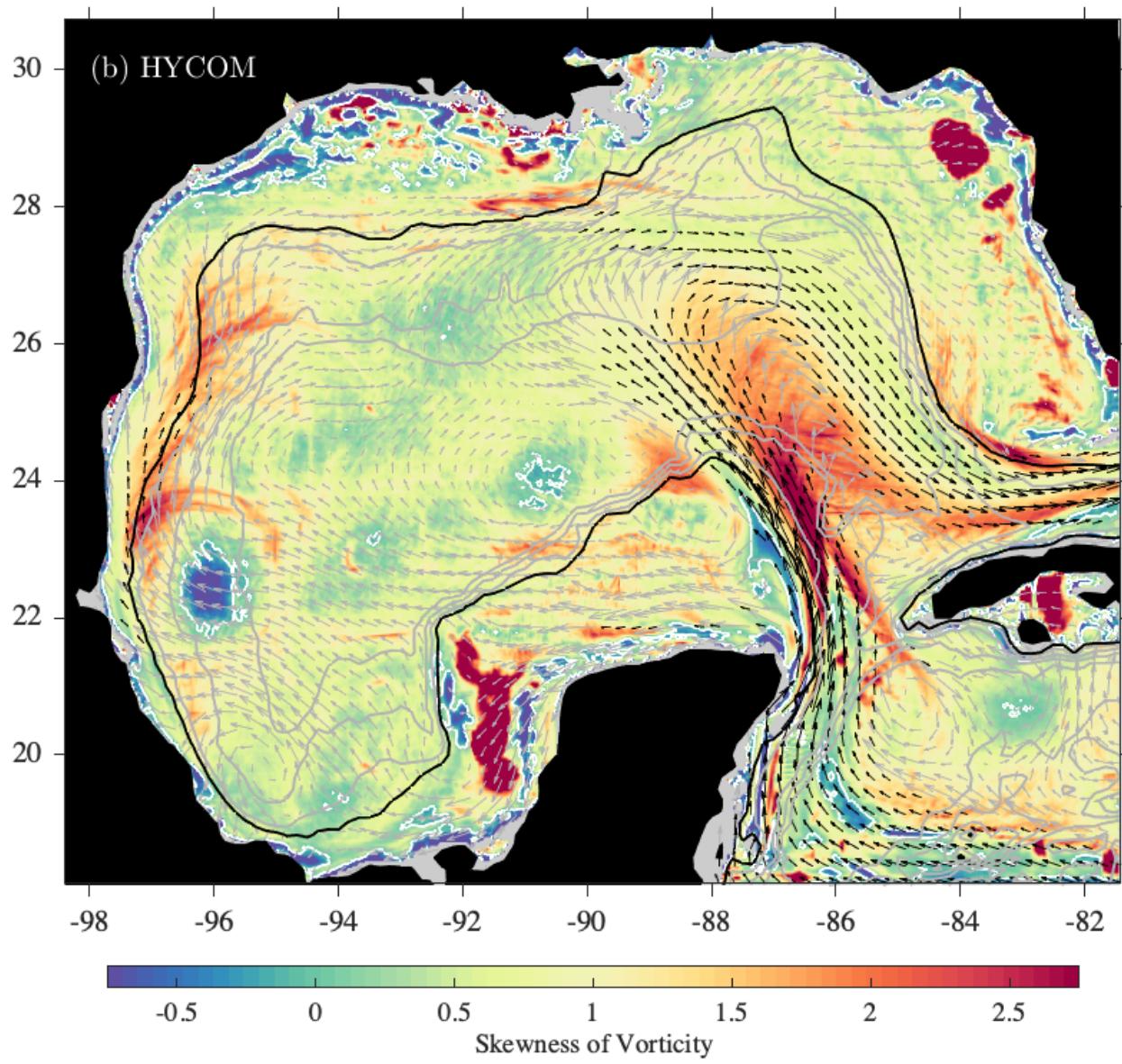
$$\sigma_\zeta^2(x, y) \equiv \frac{1}{N} \sum_{n=0}^{N-1} (\zeta_n - \bar{\zeta})^2$$

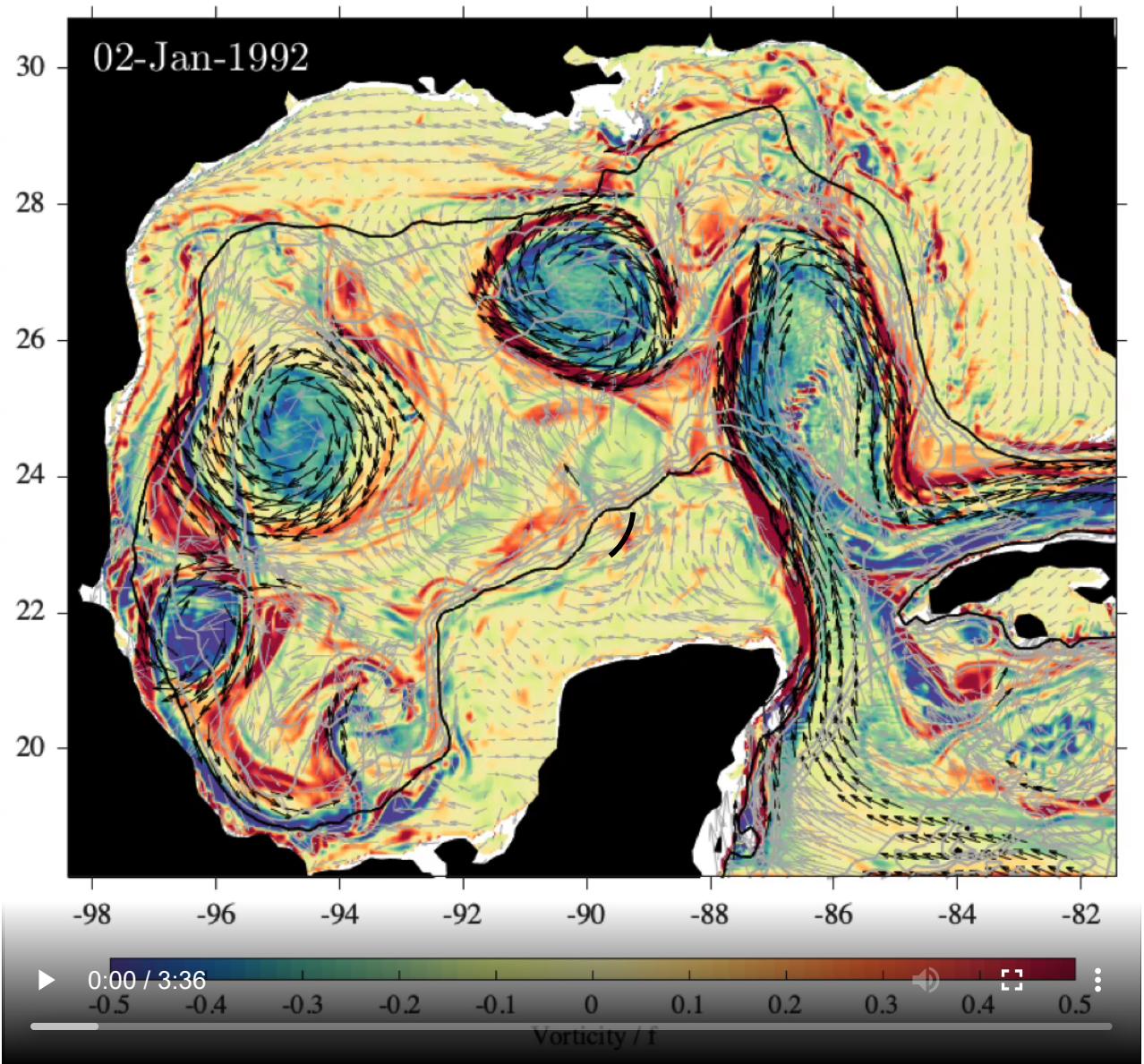
$$\bar{\gamma}_\zeta(x, y) \equiv \frac{1}{\sigma_\zeta^3} \frac{1}{N} \sum_{n=0}^{N-1} (\zeta_n - \bar{\zeta})^3.$$











# Takeaway Messages

The important message here is that time-domain statistics, while being an important tool, don't capture all the complex structure of the time-evolving turbulent ocean.

The time-domain statistics "flatten" the variability, the way a shadow flattens a three-dimensional object. They provide us with compact summaries but at the expense of compressing the rich structure.

Higher-order statistics—the skewness and kurtosis—can *sometimes* reveal features that are hidden by the lower-order statistics.

It's often a great idea to make an animation!



# Homework

1. Compute and plot the histogram of your data. In Matlab, you can do this using Matlab's `hist` or `histogram` functions. (For bivariate data, do this and the next step separately for both components.)
2. Compute the sample mean and variance.
3. Experiment with filtering your data. In Matlab, this can be done using `vfilt`. Plot the data, the filtered version, and the residual (original minus filtered) for a few choices of filter length. Are there any choices that seems to be suitable for your data?
4. Re-do the steps 1&2 involving the time-domain statistics, but using firstly the smoothed, and secondly the residual, versions of your data. How do the statistics change dependent upon the lowpass or highpass filtering? How do you interpret this?

