

# Lecture 3

## Math & Probability

## Background

ch. 1-2 of *Machine Vision* by Wesley E. Snyder & Hairong Qi

Spring 2025

16-725 (CMU RI) : BioE 2630 (Pitt)

Dr. John Galeotti



The content of these slides by John Galeotti, © 2012 - 2025 Carnegie Mellon University (CMU), was made possible in part by NIH NLM contract# HHSN276201000580P, and is licensed under a Creative Commons Attribution-NonCommercial 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc/3.0/> or send a letter to Creative Commons, 171 2nd Street, Suite 300, San Francisco, California, 94105, USA. Permissions beyond the scope of this license may be available either from CMU or by emailing [itk@galeotti.net](mailto:itk@galeotti.net).  
The most recent version of these slides may be accessed online via <http://itk.galeotti.net/>

# General notes about the book

- The book is an overview of many concepts
- Top quality design requires:
  - Reading the cited literature
  - Reading more literature
  - Experimentation & validation

# Two themes

- Consistency

- A conceptual tool implemented in many/most algorithms
- Often must fuse information from many local measurements and prior knowledge to make global conclusions about the image

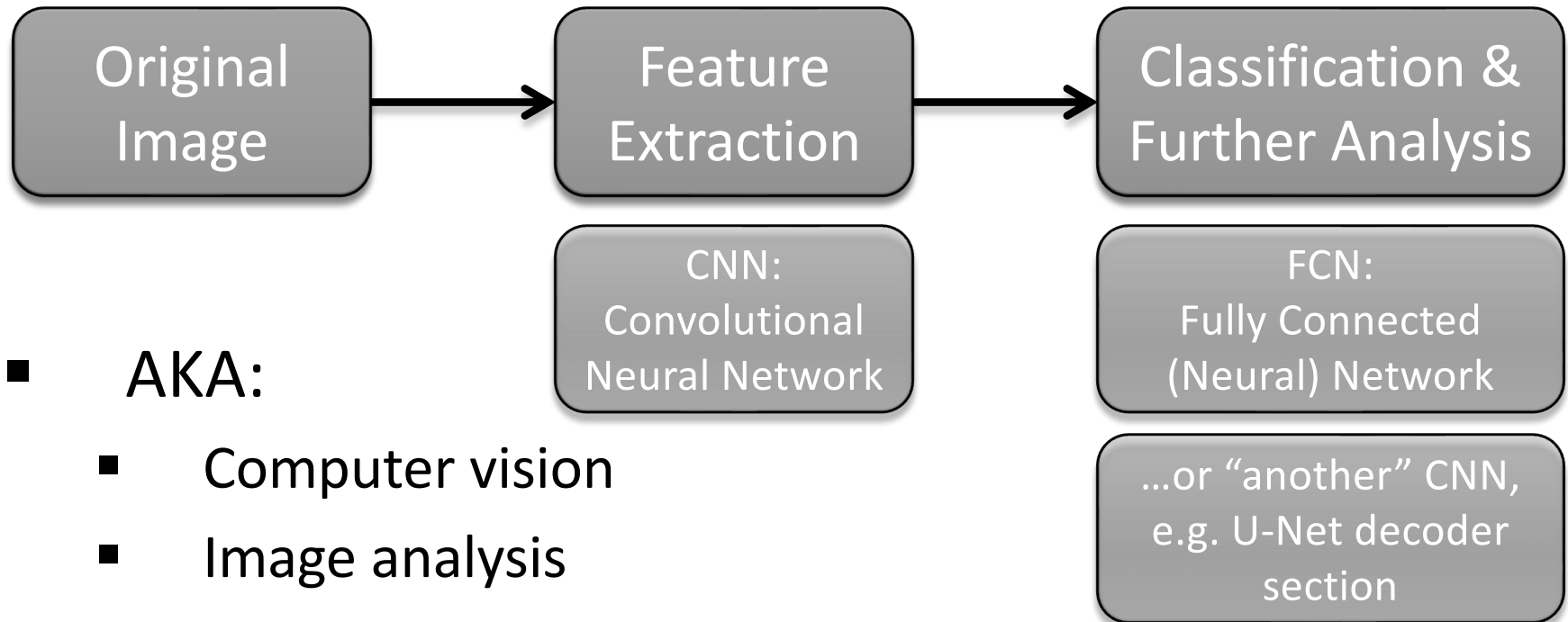
- Optimization

- Mathematical mechanism
- The “workhorse” of machine vision

# Image Processing Topics

- Enhancement
- Coding
  - Compression
- Restoration
  - “Fix” an image
  - Requires model of image degradation
- Reconstruction

# Machine Vision Topics



- AKA:

- Computer vision
- Image analysis
- Image understanding

- Pattern recognition:

1. Measurement of features

Features characterize the image, or some part of it

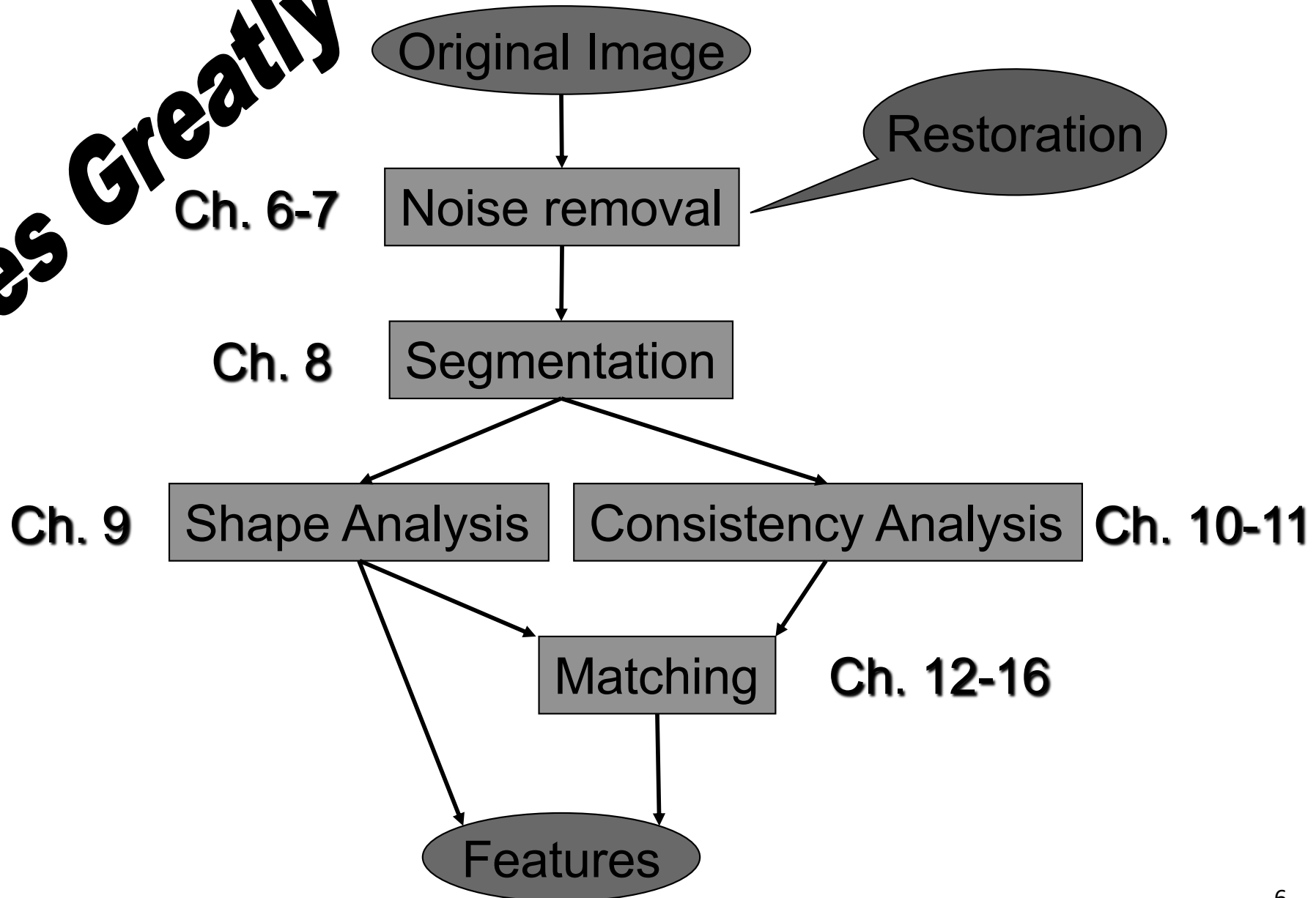
2. Pattern classification

Requires knowledge about the possible classes

Our Focus

# Feature measurement

**Varies Greatly**



# Probability

- Probability of an event  $a$  occurring:
  - $Pr(a)$
- Independence
  - $Pr(a)$  does not depend on the outcome of event  $b$ , and vice-versa
- Joint probability
  - $Pr(a,b)$  = Prob. of both  $a$  and  $b$  occurring
- Conditional probability
  - $Pr(a|b)$  = Prob. of  $a$  if we already know the outcome of event  $b$
  - Read “probability of  $a$  given  $b$ ”

# Probability for continuously-valued functions

- Probability distribution function:

$$P(x) = \Pr(z < x)$$

- Probability density function:

$$p(x) = \frac{d}{dx} P(x)$$

$$\int_{-\infty}^{\infty} p(x) dx = 1$$



# Linear algebra

$$\mathbf{v} = [x_1 \ x_2 \ x_3]^T \quad \mathbf{a}^T \mathbf{b} = \sum_i a_i b_i \quad |\mathbf{x}| = \sqrt{\mathbf{x}^T \mathbf{x}}$$

- Unit vector:  $|\mathbf{x}| = 1$
- Orthogonal vectors:  $\mathbf{x}^T \mathbf{y} = 0$
- Orthonormal: orthogonal unit vectors
- Inner product of continuous functions

$$\langle f(x), g(x) \rangle = \int_a^b f(x) g(x) dx$$

- Orthogonality & orthonormality apply here too

# Linear independence

- No one vector is a linear combination of the others
  - $\mathbf{x}_j \neq \sum a_i \mathbf{x}_i$  for any  $a_i$  across all  $i \neq j$
- Any linearly independent set of  $d$  vectors  $\{\mathbf{x}_{i=1\dots d}\}$  is a basis set that spans the space  $\mathfrak{R}^d$ 
  - Any other vector in  $\mathfrak{R}^d$  may be written as a linear combination of  $\{\mathbf{x}_i\}$
- Often convenient to use orthonormal basis sets
- Projection: if  $\mathbf{y} = \sum a_i \mathbf{x}_i$  then  $a_i = \mathbf{y}^T \mathbf{x}_i$

# Linear transforms

- = a matrix, denoted e.g.  $A$

- Quadratic form:

$$\mathbf{x}^T A \mathbf{x}$$

$$\frac{d}{d\mathbf{x}}(\mathbf{x}^T A \mathbf{x}) = (A + A^T) \mathbf{x}$$

- Positive definite:

- Applies to  $A$  if

$$\mathbf{x}^T A \mathbf{x} > 0 \quad \forall \mathbf{x} \in \mathbb{R}^d, \mathbf{x} \neq 0$$

# More derivatives

- Of a scalar function of  $\mathbf{x}$ :
  - Called the gradient
  - Really important!

$$\frac{df}{d\mathbf{x}} = \left[ \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \dots \quad \frac{\partial f}{\partial x_d} \right]^T$$

$$\frac{df}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_d} \end{bmatrix}$$

- Of a vector function of  $\mathbf{x}$ 
  - Called the Jacobian

$$\begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & \frac{\partial^2 f}{\partial x_d \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_d^2} \end{bmatrix}$$

- Hessian = matrix of 2nd derivatives of a scalar function

# Misc. linear algebra

- Derivative operators
- Eigenvalues & eigenvectors
  - Translates “most important vectors”
    - Of a linear transform (e.g., the matrix  $A$ )
  - Characteristic equation:  $A\mathbf{x} = \lambda\mathbf{x} \quad \lambda \in \mathfrak{R}$
  - $A$  maps  $\mathbf{x}$  onto itself with only a change in length
  - $\lambda$  is an eigenvalue
  - $\mathbf{x}$  is its corresponding eigenvector

# Function minimization

- Find the vector  $\mathbf{x}$  which produces a minimum of some function  $f(\mathbf{x})$

- $\mathbf{x}$  is a parameter vector
  - $f(\mathbf{x})$  is a scalar function of  $\mathbf{x}$ 
    - The “objective function”

- The minimum value of  $f$  is denoted:

$$\hat{f}(\mathbf{x}) = \min_{\mathbf{x}} f(\mathbf{x})$$

- The minimizing value of  $\mathbf{x}$  is denoted:

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x})$$

# Numerical minimization

- Gradient descent
  - The derivative points away from the minimum
  - Take small steps, each one in the “down-hill” direction
- Local vs. global minima
- Combinatorial optimization:
  - Use simulated annealing
- Image optimization:
  - Use mean field annealing
- More recent improvements to gradient descent:
  - Momentum, changing step size
- Training CNN: ADAM: an enhanced version of Stochastic Gradient Descent (SGD) w/ Momentum

# Markov models

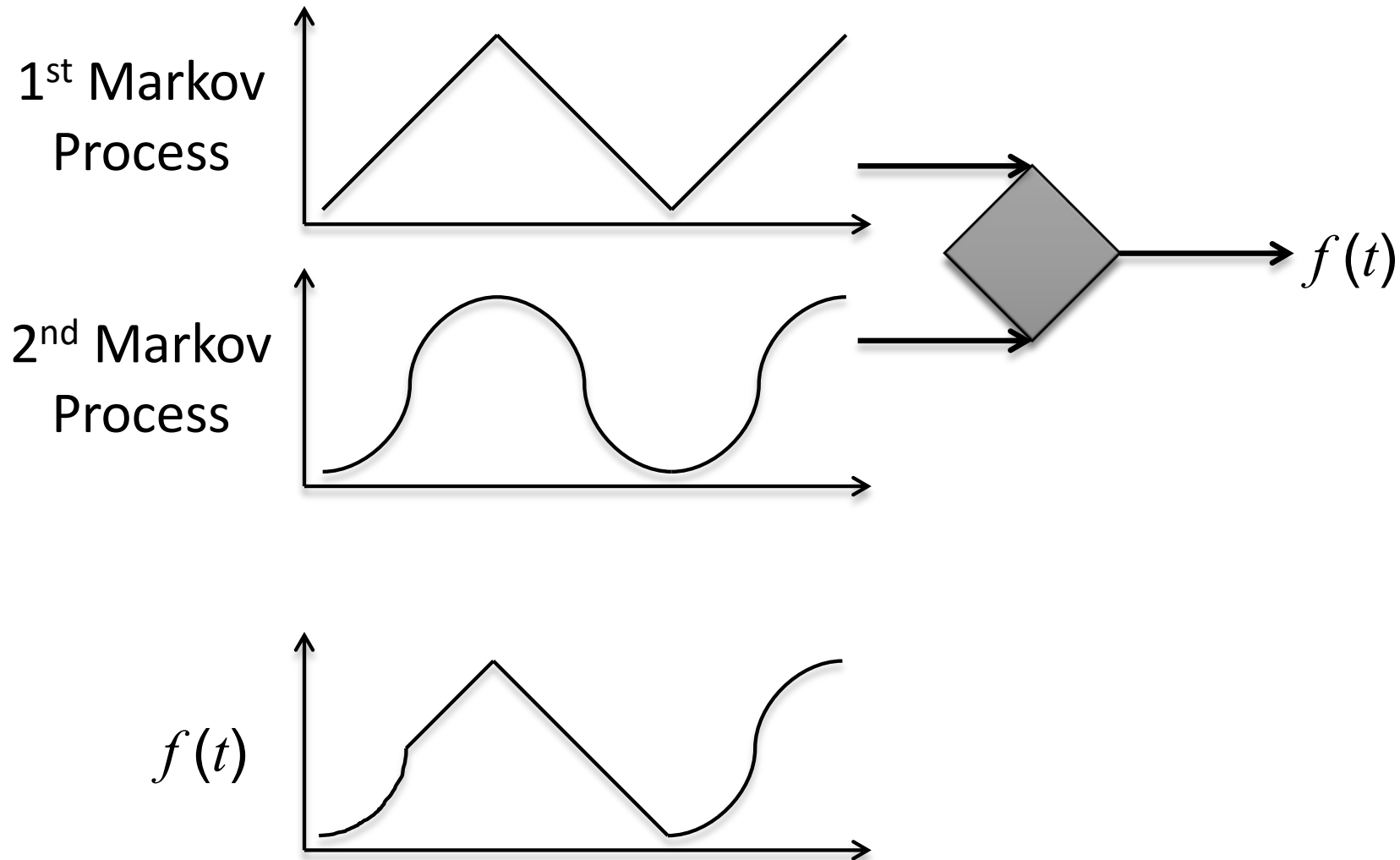
- For temporal processes:
  - The probability of something happening is dependent on a thing that just recently happened.
- For spatial processes
  - The probability of something being in a certain state is dependent on the state of something nearby.
  - Example: The value of a pixel is dependent on the values of its neighboring pixels.



# Markov chain

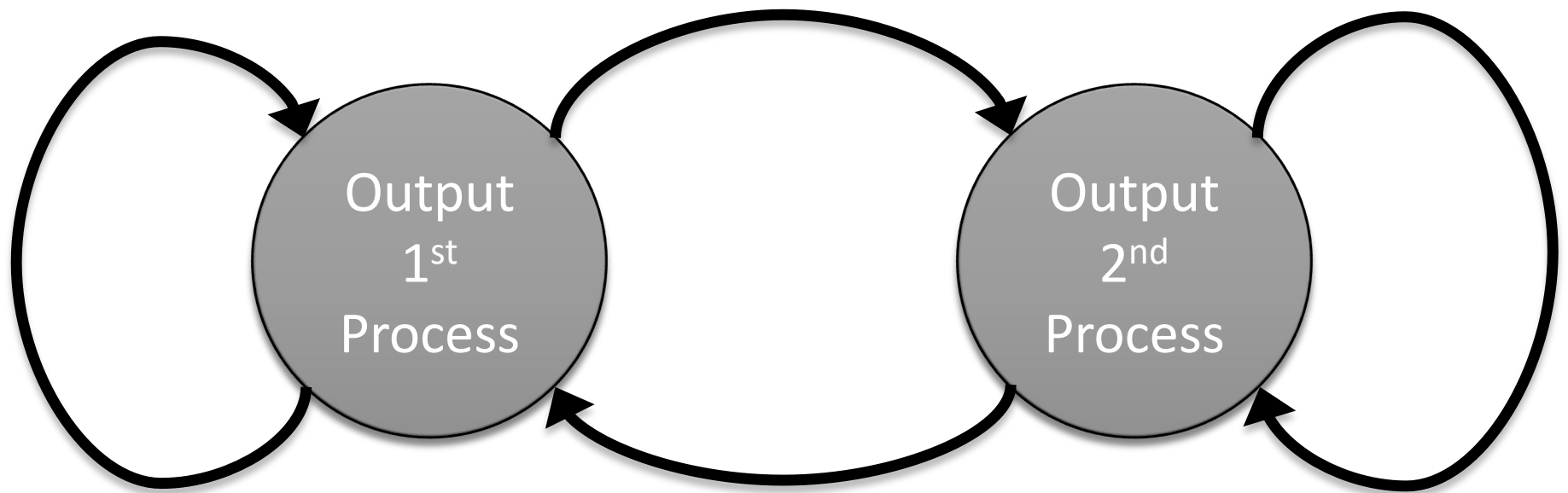
- Simplest Markov model
- Example: symbols transmitted one at a time
  - What is the probability that the next symbol will be  $w$ ?
- For a “simple” (i.e. first order) Markov chain:
  - “The probability conditioned on all of history is identical to the probability conditioned on the last symbol received.”

# Hidden Markov models (HMMs)



# HMM switching

- Governed by a finite state machine (FSM)



# The HMM Task

- Given only the output  $f(t)$ , determine:
  1. The most likely state sequence of the switching FSM
    - Use the Viterbi algorithm (much better than brute force)
    - Computational Complexity of:
      - Viterbi:  $(\# \text{ state values})^2 * (\# \text{ state changes})$
      - Brute force:  $(\# \text{ state values})^{(\# \text{ state changes})}$
  2. The parameters of each hidden Markov model
    - Use the iterative process in the book
    - Better, use someone else's debugged code that they've shared