*Week 4*
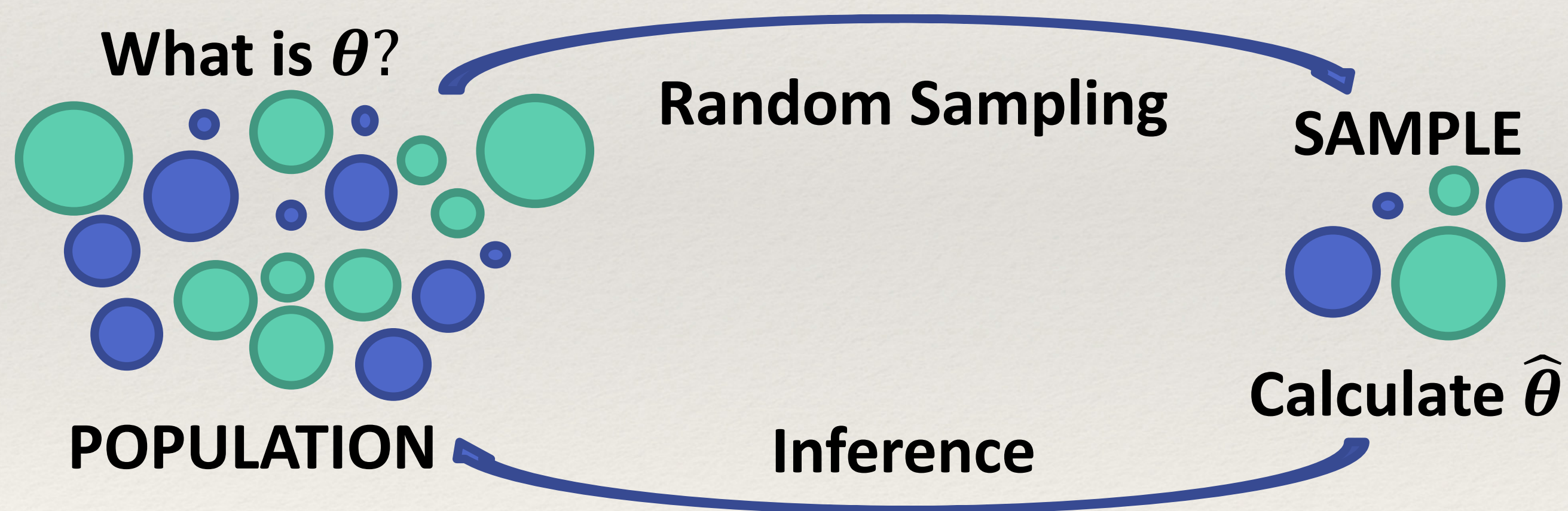
# Sampling Variability & The Central Limit Theorem

ST 314
Introduction to Statistics for Engineers

Oregon State University

# Inferential Statistics

Recall that **inferential statistics** use information from a <u>sample</u> to estimate or test characteristics from a population of interest.

**What is $\theta$?**

**Random Sampling**

**SAMPLE**

**Calculate $\hat{\theta}$**

**Inference**

**POPULATION**

How good is the sampled statistic $\hat{\theta}$ at estimating the population parameter $\theta$?

Things to consider:

❖ Method of data collection

❖ Sampling variability

❖ Sample size

# Point Estimates

As part of a quality control process for computer chips, an engineer a factory randomly samples 212 chips during a week of production to test the current rate of chips with severe defects. She finds that 27 of the chips are defective.

(a) What population is under consideration in the data set?

All chips manufactured at this factory during the week of production.

(b) What is the parameter being estimated?

$p$ = proportion of defective chips from the pop.

(c) Based on the sample what is the **point estimate** for the parameter?

$$\hat{p} = \frac{27}{212} = 0.127$$

# Sampling Variability

Suppose the study previously described was repeated by two other engineers in the same week. The following table gives the results from each of the three studies.
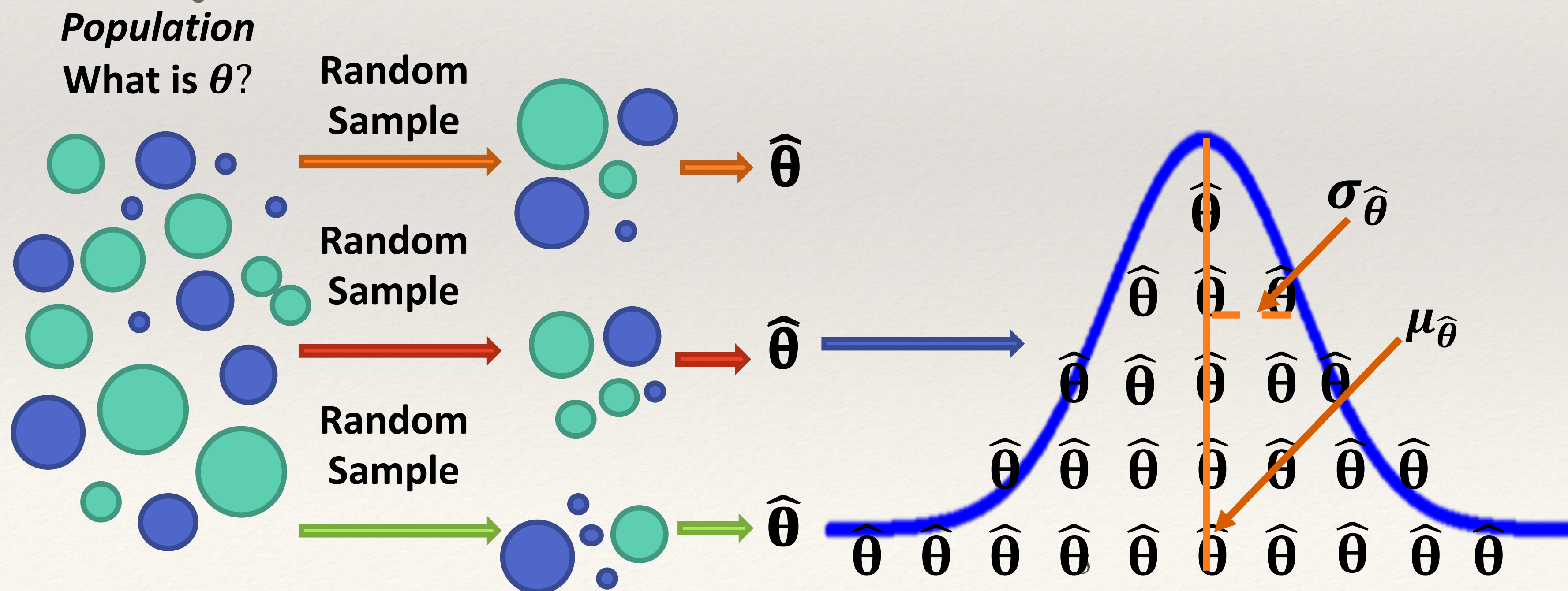
| Study | Number of Chips Sampled | Number of Defective Chips | $\hat{p}$ |
|-------|--------------------------|----------------------------|-----------|
| 1 | 212 | 27 | 0.127 |
| 2 | 212 | 19 | 0.090 |
| 3 | 212 | 23 | 0.108 |

Compare the point estimates from the three studies. What do you notice? Which point estimate is the "best"?

All point estimates are different

# Sampling Distributions

The probability distribution of a statistic, $\hat{\theta}$, is the **Sampling Distribution.** The sampling distribution defines <u>the variability of $\hat{\theta}$ and</u> <u>quantifies the chance occurrence of specific</u> values.

**Population**
**What is $\theta$?**

Random Sample → $\hat{\theta}$

Random Sample → $\hat{\theta}$

Random Sample → $\hat{\theta}$

$\sigma_{\hat{\theta}}$

$\mu_{\hat{\theta}}$

Statistics are random variables! If $N$ is the number of units in the population and $n$ is the sample size, there are $\binom{N}{n}$ possible sample combinations.

# Unbiased Estimators

❖ Common statistics, such as $\hat{p}, \bar{x},$ and $s^2$ are **unbiased.**

❖ The expected values of these estimators are equal to their parameters:

$$E(\hat{\theta}) = \theta$$

$$E(\bar{x}) = \mu \leftarrow \text{pop. mean}$$

$$E(\hat{p}) = p \leftarrow \text{pop. proportion}$$

$$E(s^2) = \sigma^2 \leftarrow \text{pop. variance}$$

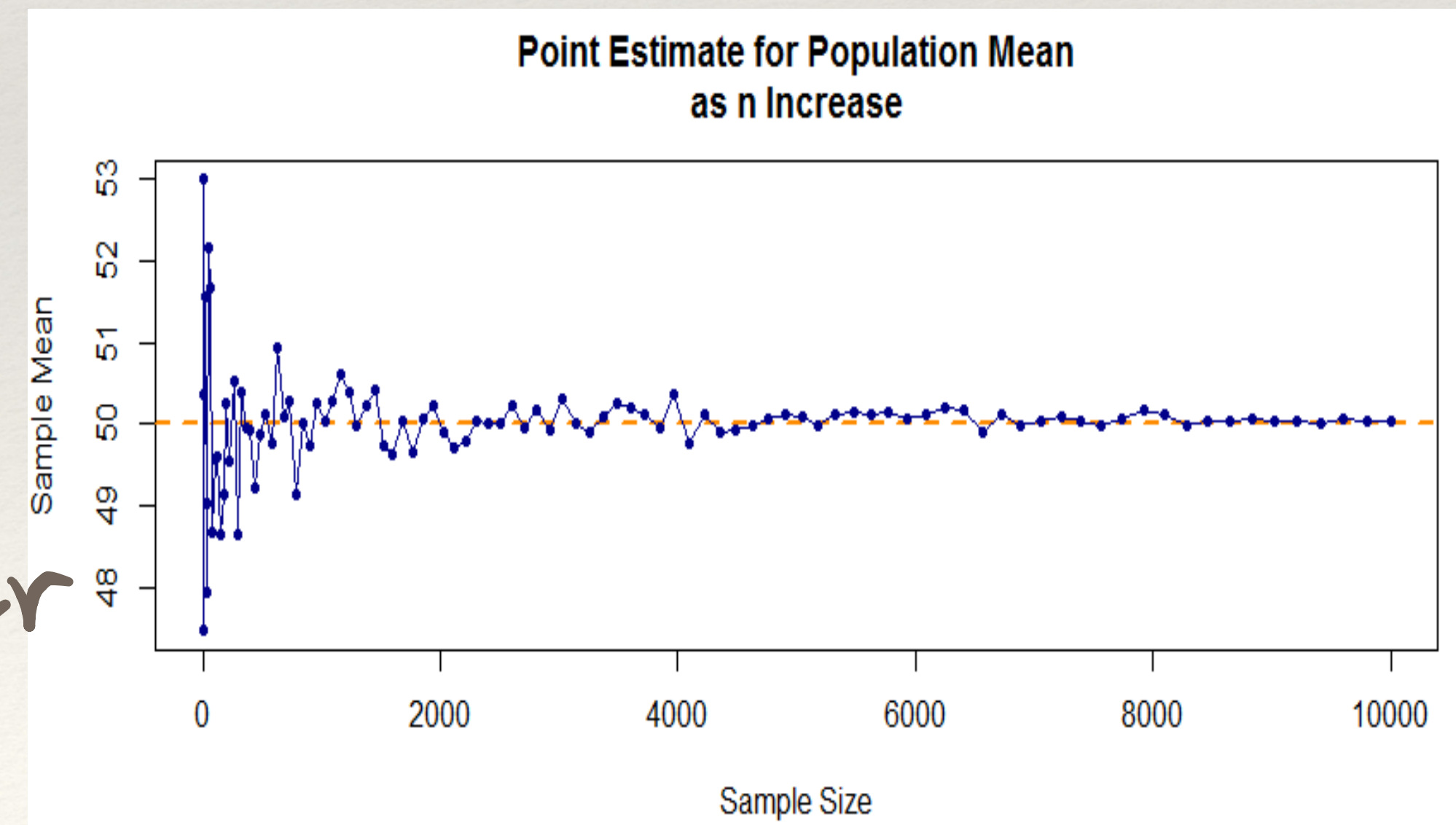$\uparrow$ sample variance

6

# Law of Large Numbers

The **Law of Large Numbers** states that as $n$ increases, the statistic will approach the true population parameter.

$$\text{As} \quad n \to \underset{\text{pop. size}}{N} \quad , \quad \hat{\theta} \to \theta$$
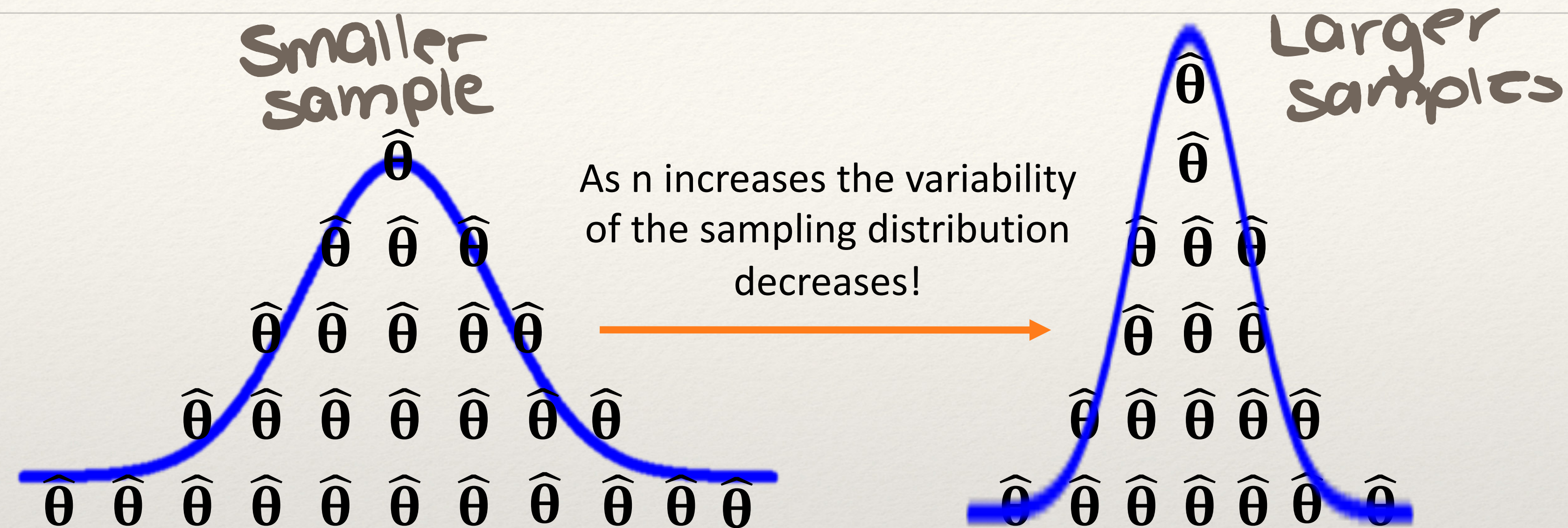
More formally, $\hat{\theta}$ converges in probability to $\theta$. This implies that $\hat{\theta}$ is a <u>consistent</u> estimator.

$$\bar{x}, \hat{p}, \text{ and } s^2 \text{ are all consistent}$$
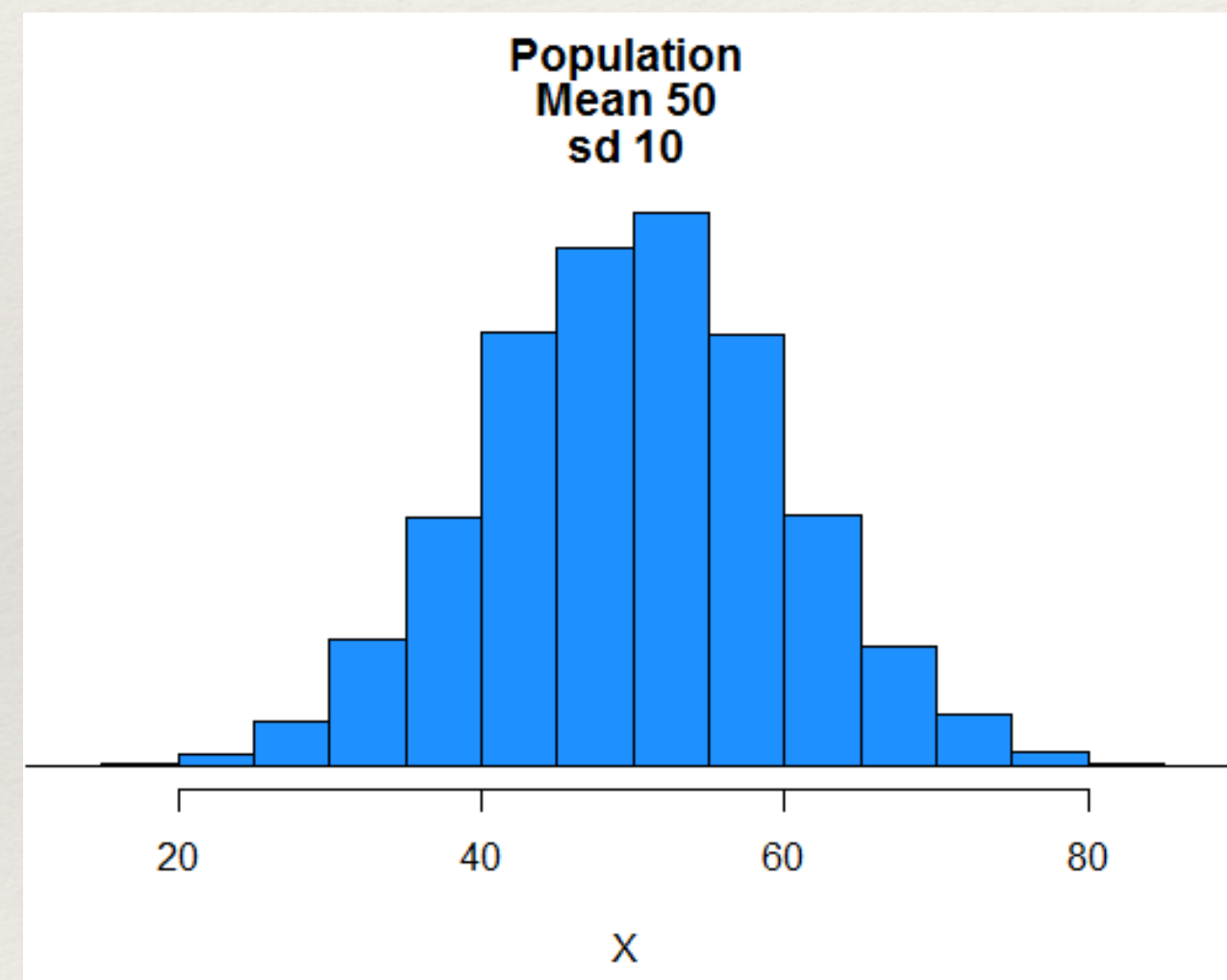
Larger sample size doesn't guarantee a closer estimate.

**! CAUTION**

Point Estimate for Population Mean as n Increase

# Sample Size & Sampling Variability

Smaller sample

$\hat{\theta}$
$\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$
$\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$
$\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$
$\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$

As n increases the variability of the sampling distribution decreases!

Larger samples

$\hat{\theta}$
$\hat{\theta}$
$\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$
$\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$
$\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$
$\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$ $\hat{\theta}$

The variability of the sampling distribution of $\hat{\theta}$ is referred to as the **standard error** _____, denoted by $SE_{\hat{\theta}}$ or $\sigma_{\hat{\theta}}$. The standard error describes the typical error or uncertainty of the statistic. It is the standard deviation of the statistic.

8

# Sample Size & Sampling Variability

A random variable X is simulated from a normal distribution with population parameters $\mu_x$ and $\sigma_x$.



As *n* increases, the variability of the sampling distribution decreases .

# Distributions in inference
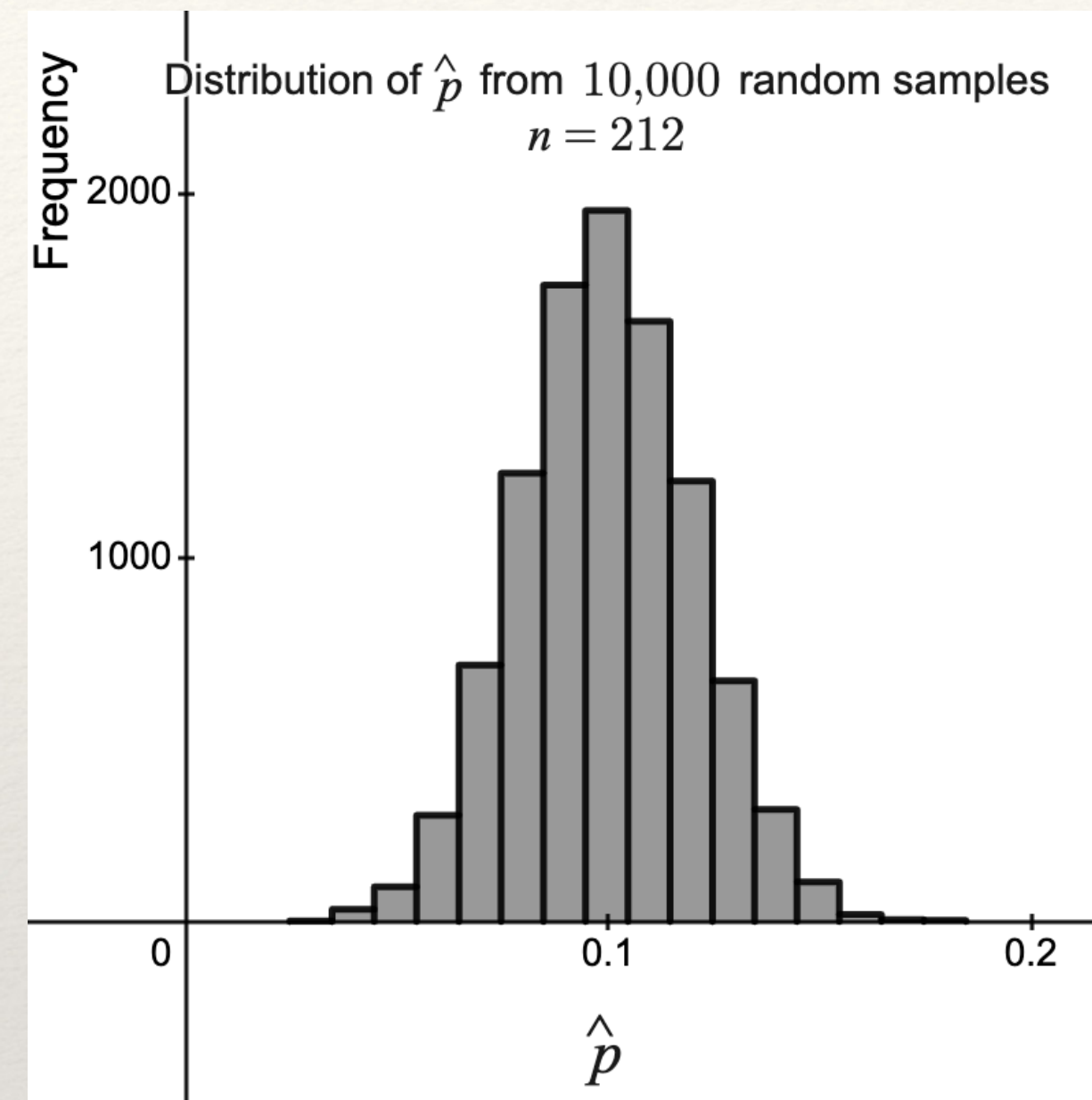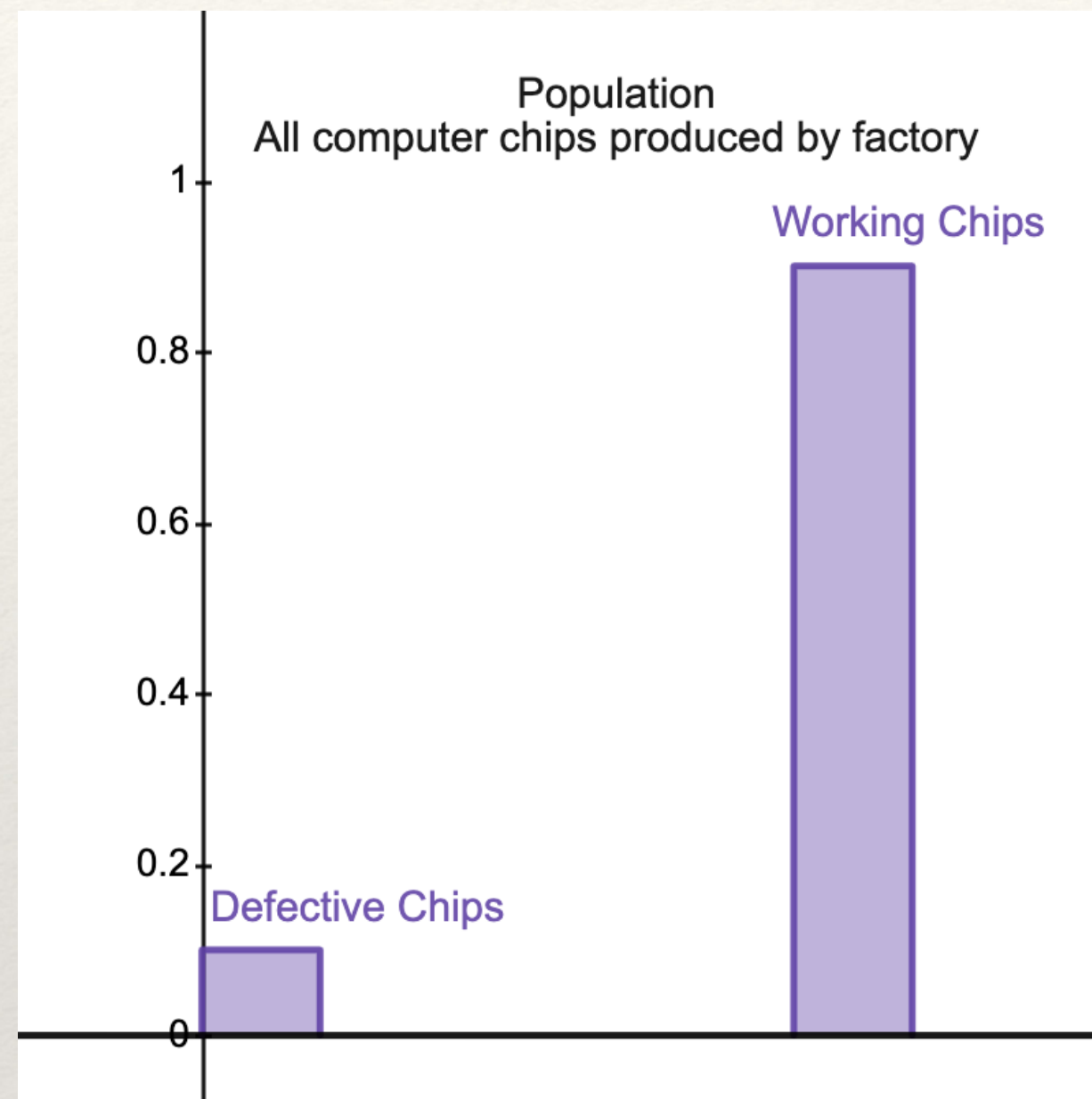


| Population Distribution | Sampled Distribution | Sampling Distribution |
|---|---|---|
| Distribution of the entire collection of observations | Distribution of n obs. obtained from a single sample | Distribution of a sampled statistic from repeated samples of size n from population |

# Sampling Distribution of the Sample Proportion



Population
All computer chips produced by factory

Working Chips

Defective Chips



Distribution of $\hat{p}$ from $10{,}000$ random samples
$n = 212$

Frequency

$\hat{p}$

If $n$ is sufficiently large, then the **Central Limit Theorem** states the sampling distribution of the statistic $\hat{p}$ is:

$$\hat{p} \sim N\left(p, \sqrt{\frac{p(1-p)}{n}}\right)$$

# Sampling Distribution of the Sample Mean



**Population Data**



Distribution of $\bar{x}$ from 10,000 random samples n=35

If $n$ is sufficiently large, then the **Central Limit Theorem** states the sampling distribution of the statistic $\bar{x}$ is:

$$\bar{x} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

population mean

pop. standard dev.

# Sufficiently Large Samples

$n\hat{p} \geq 10$

$n(1-\hat{p}) \geq 10$

| **Normal Data** $n \geq 1$ | **Symmetric Data** $n \geq 12$ | **Skewed Data** $n \geq 30$ | **Bernoulli Data** $n \geq \dfrac{10}{p}$ and $\dfrac{10}{1-p}$ |
|---|---|---|---|

To apply CLT, sample size must be sufficiently large



Whether $n$ is **sufficiently large enough** is determined by the population shape, excluding binary data, $n \geq 30$ is the rule of thumb.

# Sampling Distributions

Sampling distributions are never observed, but we keep them in mind!