# NIH Collaboratory Rethinking Clinical Trials®

Health Care Systems Research Collaboratory

#### Living Textbook Grand Rounds Series

#### Demystifying Biostatistical Concepts for Embedded Pragmatic Clinical Trials

#### June 19, 2020

Elizabeth L. Turner, PhD, Duke University Patrick J. Heagerty, PhD, University of Washington David M. Murray, PhD, National Institutes of Health

For the NIH Collaboratory Coordinating Center Biostatistics and Study Design Core Working Group

#### Overview

- Focus of this talk: demystifying design-related issues for embedded pragmatic clinical trials (ePCTs)
- Context: NIH Collaboratory–funded studies
- Three kinds of randomized trials
  - Randomized controlled trial (RCT)
  - Cluster randomized trial (CRT)
    - Parallel vs stepped-wedge
  - Individually randomized group treatment (IRGT) trial
- How to select amongst these designs?
- Other brief topics: clustering, power, and analytical issues



## In the Living Textbook

#### DESIGN

#### EXPERIMENTAL DESIGNS AND RANDOMIZATION SCHEMES

- 1 Introduction
- 2 Statistical Design Considerations
- **3** Cluster Randomized Trials
- 4 Randomization Methods
- 5 Choosing Between Cluster and Individual Randomization
- 6 Alternative Cluster Randomized Designs
- 7 Concealment and Blinding
- 8 Designing to Avoid Identification Bias
- 9 Additional Resources

#### ANALYSIS PLAN

- 1 Introduction
- 2 Intraclass Correlation
- 3 Unequal Cluster Sizes
- 4 Accounting for Residual Confounding in the Analysis
- 5 Missing Data and Intention-to-Treat Analyses
- 6 EHR Data Extraction
- 7 Unanticipated Changes
- <sup>8</sup> Case Study: STOP CRC Trial

#### NIH Collaboratory ePCT: SPOT

SUICIDE PREVENTION OUTREACH TRIAL

- Suicide Prevention Outreach Trial (SPOT)
- Approximately 16,000 patients across 4 clinical sites
- Three-arm RCT to evaluate 2 individual-level interventions vs usual care
- Interventions
  - Skills training program
  - Care management program
- Intervention contact mostly though EHR
  - Low risk of "contamination"
  - Individual-level randomization appropriate
- Unit of randomization: patient

Simon GE et al. *Trials*. 2016;17(1):452.

#### NIH Collaboratory ePCT: STOP CRC



- Strategies and Opportunities to Stop Colorectal Cancer in Priority Populations (STOP CRC)
- 40,000+ patients across 26 clinical sites
- Intervention
  - Health system—based program to improve CRC screening rates
  - Applied to clinical site  $\rightarrow$  cluster randomization
- Unit of randomization: clinical site
- Two-arm cluster randomized trial (CRT)
  - Also referred to as a group-randomized or community randomized trial

Coronado GD et al. Contemp Clin Trials. 2014;38(2):344-349.

## Reasons to Randomize Clusters Instead of Individuals

- Intervention targets health care units rather than individuals
  - STOP CRC: clinic-based intervention to improve screening
- Intervention targeted at individual at risk of contamination
  - Intervention adopted by members of control arm
  - For example, physicians randomized to new educational program may share knowledge with control-arm physicians in their practice
  - Contamination reduces the observed treatment effect
- Logistically easier to implement intervention by cluster

#### **STOP CRC Cluster Randomization**



**Level 2**: Randomization at the level of the clinic (ie, cluster)



**Level 1**: Individual-level outcomes nested within clinics

#### **STOP CRC Cluster Randomization**



- Level 1: Individual-level outcomes nested within clinics
- Individual-level outcomes within same clinic expected to be correlated (ie, to *cluster*)

#### **STOP CRC Cluster Randomization**



- Level 1: Individual-level outcomes nested within clinics
- Individual-level outcomes within same clinic expected to be correlated (ie, to *cluster*)
- Reduces power to detect treatment effect if same sample size used as under individual randomization

#### **Understanding Outcome Clustering**

- Consider 10 control-arm clinics (ie, clusters)
- Each with 5 age-eligible patients: ie, who are not up to date with colorectal cancer (CRC) screening
- Binary outcome: refused screening (Y/N)

# Understanding Outcome Clustering: <u>Complete</u> Clustering



Screened
Not screened

# Understanding Outcome Clustering: <u>Complete</u> Clustering



>1 participant/clinic gives no more information than a single participant/clinic since every participant in a given clinic has the same outcome

# Understanding Outcome Clustering: <u>No</u> Clustering



Not screened

# Understanding Outcome Clustering: <u>No</u> Clustering



Not screened

20% uptake of CRC screening in each clinic No structure by clinic; more like a random sample of eligible participants

# Understanding Outcome Clustering: <u>Some</u> Clustering



Not screened

# Understanding Outcome Clustering: <u>Some</u> Clustering



A more typical situation: proportion screened ranges from 0% - 80%

# Measure of Outcome Clustering: Intraclass Correlation Coefficient (ICC)

- Needed for study planning and power
- Most commonly used measure of clustering
- Ranges: 0-1; 0 = no clustering; 1 = complete clustering
- Typically < 0.2; commonly around 0.01 to 0.05
- Between-cluster outcome variance vs total outcome variance

# Measure of Outcome Clustering: Intraclass Correlation Coefficient (ICC)

- Needed for study planning and power
- Most commonly used measure of clustering
- Ranges: 0-1; 0 = no clustering; 1 = complete clustering
- Typically < 0.2; commonly around 0.01 to 0.05
- Between-cluster outcome variance vs total outcome variance

ICC for continuous outcomes:

$$\Gamma = \frac{S_B^2}{S_B^2 + S_W^2} = \frac{S_B^2}{S_{Total}^2}$$

Involves both between-cluster and within-cluster variance



	DESIGN		He
ANAL	YSIS PLAN		
1	Introduction		
2	Intraclass Correlation		
3	Unequal Cluster Sizes	J	
4	Accounting for Residual Confounding in the Analysis		
5	Missing Data and Intention-to-Treat Analyses		
6	EHR Data Extraction		
7	Unanticipated Changes		
8	Case Study: STOP CRC Trial		

#### 14 NIH Collaboratory<sub>Rethinking Clinical Trials</sub> ealth Care Systems Research Collaboratory

#### Intraclass Correlation Coefficient Cheat Sheet

This document provides an introductory description of the intraclass correlation coefficient (ICC), a descriptive statistic that is Inis document provides an introductory description of the intractass correlation coefficient (ILL), a descriptive statistic that is important for the design and analysis of cluster randomized trials. In a <u>cluster randomized trial</u>, instead of being randomized intractant descriptions the units of conductivation is a cluster such as a second of contribution being can also be provided endowing and analysis of the units of conductivation is a cluster andomized trial. Important for the design and analysis of cluster-randomized trials. In a <u>cluster randomized trial</u>, instead of being randomized by individual participant, the unit of randomization is a cluster, such as a group of participants being seen at a hospital, clinic, the component of the other sector and the participant of the provided sector individual land. or primary care practice, although the outcomes may still be measured at an individual level.

The intractass complation coefficient (ICC) is a descriptive statistic that describes the extent to which outcomes 1) within each interest an illuster has complete the second of the end The intractass correlation coefficient (ICC) is a descriptive statistic that describes the extent to which outcomes 1) within each duster are likely to be similar or 2) between different clusters are likely to be *different* from each other, relative to outcomes from other division. The ICC is an important tool for divider endowing descention trials because this value halos determined in the division. cluster are likely to be similar or 2) between different clusters are likely to be *sigrerent* from each other, relative to outcomes from other clusters. The ICC is an important tool for cluster randomized pragmatic trials because this value helps determine the component of the signal and the from other clusters. The ICC is an important tool for cluster-randomized pragmatic trials because this value nelps determine the sample size needed to detect a treatment effect. Although it ranges from 0 to 1 theoretically, the ICC for most pragmatic clusters endowland birds in the same of a composite second on the 0 ne tine sample size measure to weters a treatment eners. Autoorgn is ranges in duster-randomized trials is typically (0.2; commonly around 0.01 to 0.05.

In duster-randomized trials where groups of individuals are randomized to treatment arms, when outcomes within dusters biother and and unless the meaning of individuals of endowing a force duster is write different them explicit endow unless the second In cluster-randomized trials where groups of individuals are randomized to treatment arms, when outcomes writin custers are highly correlated and when the magnitude of outcomes across clusters is quite different, then participants within the during set links to how circlar outcomes and the VC will be been when this is the case, the data from one member of the are highly correlated and when the magnitude of outcomes across clusters is quite different, then participants within the duster are likely to have similar outcomes and the ICC will be large. When this is the case, the data from one member of the duster are likely to have similar outcomes and the ICC will be large. When this is the case, the data from one member of the outcome duster are likely to have similar outcomes and the ICC will be large. When this is the case, the data from one member of the outcome duster are likely to have similar outcomes and the outcome are instituted. Uncome the different computer size is duster to the outcome of the outcom cluster are likely to have similar outcomes and the ICC will be large. When this is the case, the data from one member of the cluster provides almost as much information as if all of the members are included. Hence, the effective sample size is closer to be summarized clusters of expended to be active scenario size of study constringents. vusiver provides annuas cas much information as n an or the memores are included. It to the number of clusters as opposed to the entire sample size of study participants.

To demonstrate why this is relevant, let's consider two examples:

- 1 In a dietary intake study, the data from several members in a dietary intake study, die wata north schellar and would of the same family would likely be very similar and would of the same failing would likely be very annual and woo differ from that of other families. Hence there may be little gain from sampling more than one member. On the other hand, if a cluster is an entire city and subjects within the city are randomly sampled, one might expect relatively little similarity from subject to subject relative to the rest of the sample. In this case, each individual subject would likely contribute "independent" information.
- 2 Suppose we have 6 providers, each with 3 eligible participants for a pragmatic cluster-randomized trial. In this hypothetical case, the outcome is patient satisfaction rated on a scale from 1 to 10 with an outcome distribution as shown in Figure 1. One might expect that patients seen by a specific provider will have more similar levels of satisfaction to each other than to patients from other providers and that some providers will have consistently interesting the patient satisfaction (e.g. provider high patient satisfaction (e.g. provider 2) whereas others will have consistently low patient satisfaction (e.g. provider 2) whereas others will have consistent be be applied on the source of the sourc ingri patient satisfaction (e.g., provider 4) whereas ounces with nave consistently now patient satisfaction (e.g., provider 1). This is an example of how outcomes within each cluster are likely to be similar. Thus, the ICC is high, and adding cluster are to be a first stream does not encode an each dustriant information. It is an example or now outcomes when each case are need to individuals to the cluster does not provide much additional information.



Accounting for Clustering Requires Larger Sample for Adequate Power

- Power and detectable difference is affected by...
  - Strength of the clustering effect (eg, size of ICC)
  - Number of clusters
  - Number of patients per cluster











Example: CRT with smaller ICC=0.01 at at fixed alpha & power



#### Impact of increasing # clusters/groups

Example: CRT with even smaller ICC=0.001 at fixed alpha & power



#### Accounting for Clustering in Design

- Power and sample size for CRT
  - Account for anticipated clustering
  - Inflate RCT sample size
  - Work with statistician to do correctly
- Use ICC for outcome
  - ICC often 0.01-0.05
  - STOP CRC: ICC = 0.03 for primary outcome
  - Depends on outcome and study characteristics
  - Different outcome = different ICC, even in same CRT

#### **Estimating ICC to Plan Study**

- How to get good estimate of ICC for a particular outcome?
  - Depends on outcome and study characteristics
  - CONSORT statement recommends ICC reported
  - Look at other articles with similar settings
  - Use available EHR data
- Be cautious when using pilot data from small study
  - ICC might have a wide confidence interval

#### NIH Collaboratory ePCT: LIRE



- Lumbar Imaging with Reporting of Epidemiology (LIRE)
- Goal: reduce unnecessary spine interventions by providing info on prevalence of normal findings
- Patients of 1700 PCPs across 100 clinics
- Clinic-level intervention  $\rightarrow$  cluster randomization
- Unit of randomization: clinic
- Pragmatic trial
  - All clinics will eventually receive intervention
  - Stepped-wedge CRT

Jarvik JG et al. Contemp Clin Trials. 2015;45(Pt B):157-163.

#### NIH Collaboratory ePCT: LIRE

Exposed to LIRE intervention

Unexposed to LIRE intervention



Source: Jarvik JG et al. Contemp Clin Trials. 2015;45(Pt B):157-163.



Parallel

Stepped-wedge



In complete designs, measurements are taken from every cluster at every time point. In incomplete designs, some clusters do not provide measurements at all time points.

Examples with 8 clusters: 1-year intervention



Based on: Hemming K, Lilford R, Girling AJ. 2015. Stepped-wedge cluster randomised controlled trials: a generic framework including parallel and multiple-level designs. Stat Med. 34:181-196. doi:10.1002/sim.6325. PMID: 25346484

Examples with 8 clusters: 1-year intervention



Based on: Hemming K, Lilford R, Girling AJ. 2015. Stepped-wedge cluster randomised controlled trials: a generic framework including parallel and multiple-level designs. Stat Med. 34:181-196. doi:10.1002/sim.6325. PMID: 25346484

Examples with 8 clusters: 1-year intervention

Control period

Intervention period



Incomplete steppedwedge design



Based on: Hemming K, Lilford R, Girling AJ. 2015. Stepped-wedge cluster randomised controlled trials: a generic framework including parallel and multiple-level designs. Stat Med. 34:181-196. doi:10.1002/sim.6325. PMID: 25346484

Examples with 8 clusters: 1-year intervention

Control period

Intervention period



Based on: Hemming K, Lilford R, Girling AJ. 2015. Stepped-wedge cluster randomised controlled trials: a generic framework including parallel and multiple-level designs. Stat Med. 34:181-196. doi:10.1002/sim.6325. PMID: 25346484

Examples with 8 clusters: 1-year intervention



Estimated (primarily) using between- cluster ie, **vertical** information





Complete SW design

**Control** period

Estimated (primarily) using between- cluster ie, **vertical** information





Complete SW design

Control period

Estimated (primarily) using between- cluster ie, **vertical** information



Estimated using both vertical & horizontal (ie, within-cluster) information



Complete SW design

Control period

Estimated (primarily) using between- cluster ie, **vertical** information



Estimated using both vertical & horizontal (ie, within-cluster) information



Complete SW design

Control period

Estimated (primarily) using between- cluster ie, **vertical** information



Estimated using both vertical & horizontal (ie, within-cluster) information



Complete SW design

Control period

Estimated (primarily) using between- cluster ie, **vertical** information

![](_page_43_Figure_2.jpeg)

Estimated using both vertical & horizontal (ie, within-cluster) information

![](_page_43_Figure_4.jpeg)

Complete SW design

Control period Intervention period

## Choosing the Right Type of CRT

- Arguments <u>for</u> stepped-wedge CRT:
  - Cannot immediately implement intervention in 1/2 clusters
  - Pragmatic research: eventually implement in all clusters
  - Have few clusters and might gain power

## Choosing the Right Type of CRT

#### Arguments <u>for</u> stepped-wedge CRT:

- Cannot immediately implement intervention in 1/2 clusters
- Pragmatic research: eventually implement in all clusters
- Have few clusters and might gain power
- Arguments <u>against</u> stepped-wedge CRT:
  - Risk confounding treatment effect with time effect
  - Risk of interruption or external events that could affect the outcome (eg, a pandemic!)

#### **Recommendations for CRT Design**

- Use a parallel CRT design if you can
- If stepped-wedge, plan for time effects in design & analysis
- Work with statistician to account for clustering in design and analysis of both designs

Is there a strong rationale for randomizing groups/clusters rather than individuals to study conditions?

Is there a strong rationale for randomizing groups/clusters rather than individuals to study conditions?

![](_page_48_Figure_2.jpeg)

Is there a strong rationale for randomizing groups/clusters rather than individuals to study conditions?

Yes

Examples with clinic/health-system-level interventions:

- STOP CRC colorectal cancer screening CRT
- LIRE lumbar imaging trial SW-CRT

Is there a strong rationale for randomizing groups/clusters rather than individuals to study conditions?

> Is there a strong rationale for rolling out the intervention to all clusters before the end of the trial?

Yes

Is there a strong rationale for randomizing groups/clusters rather than individuals to study conditions?

> Is there a strong rationale for rolling out the intervention to all clusters before the end of the trial?

#### STOP CRC colorectal cancer screening CRT

![](_page_51_Figure_4.jpeg)

Yes

Is there a strong rationale for randomizing groups/clusters rather than individuals to study conditions?

> Is there a strong rationale for rolling out the intervention to all clusters before the end of the trial?

> > Yes

SW-CRT

Yes

LIRE lumbar imaging SW-CRT

No

Is there a strong rationale for randomizing groups/clusters rather than individuals to study conditions?

No

Is there a strong rationale for randomizing groups/clusters rather than individuals to study conditions?

Examples with individual-level randomization:

- SPOT suicide prevention RCT
- OPTIMUM mindfulness for back-pain RCT

Is there a strong rationale for randomizing groups/clusters rather than individuals to study conditions?

Do participants receive their treatment in a group format or from a shared interventionist?

No

Is there a strong rationale for randomizing groups/clusters rather than individuals to study conditions?

Do participants receive their treatment in a group format or from a shared interventionist?

No

![](_page_56_Figure_3.jpeg)

SPOT suicide prevention RCT Intervention is targeted at the individual

![](_page_57_Figure_0.jpeg)

![](_page_58_Figure_0.jpeg)

OPTIMUM mindfulness for back-pain RCT Intervention is group-based

#### NIH Collaboratory ePCT: OPTIMUM

- OPTIMUM: optimizing pain treatment in medical settings using group-based mindfulness
- ~450 patients across 3 clinical sites
- Two-arm RCT
  - Intervention vs usual care
- Unit of randomization: individual
- Group-based intervention
  - Clustering of outcomes in intervention arm
  - Must be accounted for in both design and analysis
- "Individually randomized group treatment (IRGT) trial"

![](_page_59_Picture_10.jpeg)

![](_page_60_Figure_1.jpeg)

See Figure: Murray DM, Taljaard M, Turner EL, George SM, Ann Rev Pub Health 2020. 41:1-19

![](_page_61_Figure_0.jpeg)

See Figure: Murray DM, Taljaard M, Turner EL, George SM, Ann Rev Pub Health 2020. 41:1-19

# 66 Important Things to Know

- Question drives design, design drives analysis
- Randomization
  - Individual-level preferred for statistical reasons
  - But cluster randomization often needed
- Account for clustering in design and analysis of:
  - CRT
  - IRGT trial
- Good design is difficult but critical
  - Need input from diverse team, including statistician
  - Analysis may not be able to overcome design flaws

# Important Things to Do

- Focus on the research question
- Select design features with analysis in mind
- Collaborate early with a statistician
- Choose individual randomization, but only if possible
- Weigh statistical choices vs implementation challenges
- Write and publish a protocol paper

![](_page_64_Figure_0.jpeg)

## In the Living Textbook

 $\mathbf{O}$ 

#### DESIGN

#### EXPERIMENTAL DESIGNS AND RANDOMIZATION SCHEMES

- 1 Introduction
- 2 Statistical Design Considerations
- **3** Cluster Randomized Trials
- 4 Randomization Methods
- 5 Choosing Between Cluster and Individual Randomization
- 6 Alternative Cluster Randomized Designs
- 7 Concealment and Blinding
- 8 Designing to Avoid Identification Bias
- 9 Additional Resources

#### ANALYSIS PLAN

- 1 Introduction
- 2 Intraclass Correlation
- 3 Unequal Cluster Sizes
- 4 Accounting for Residual Confounding in the Analysis
- 5 Missing Data and Intention-to-Treat Analyses
- 6 EHR Data Extraction
- 7 Unanticipated Changes
- <sup>8</sup> Case Study: STOP CRC Trial

#### Summary

- Focus of this talk: demystifying design-related issues for embedded pragmatic clinical trials (ePCTs)
- Context: NIH Collaboratory–funded studies
- Three kinds of randomized trials
  - Randomized controlled trial (RCT)
  - Cluster randomized trial (CRT)
    - Parallel vs stepped-wedge
  - Individually randomized group treatment (IRGT) trial
- How to select amongst these designs?
- Other brief topics: clustering, power, and analytical issues

#### **Design and Analysis Methods**

- Turner EL et al. Review of recent methodological developments in group-randomized trials: part 1design. <u>Am J Public Health</u>. 2017;107(6):907-915.
- Turner EL et al. Review of recent methodological developments in group-randomized trials: part 2analysis. <u>Am J Public Health. 2017;107(7):1078-</u> <u>1086.</u>
- Murray DM et al. Essential ingredients and innovations in the design and analysis of grouprandomized trials. <u>Annu Rev Public Health.</u> <u>2020;41:1-19.</u>
- Li F et al. Mixed-effects models for the design and analysis of stepped wedge cluster randomized trials: an overview. *Stat Methods Med Res.* In press.
- Hemming et al. The Shiny CRT Calculator: Power and Sample size for Cluster Randomised Trials. <u>https://clusterrcts.shinyapps.io/rshinyapp/</u>

![](_page_66_Picture_6.jpeg)

#### **NIH Resources**

- Pragmatic and Group-Randomized Trials in Public Health and Medicine
  - https://prevention.nih.gov/grt
  - 7-part online course on GRTs and IRGTs
- Mind the Gap Webinars
  - <u>https://prevention.nih.gov/education-training/methods-mind-gap</u>
    - Analytic methods for SW-GRTs (Fan Li, July 14, 2020)
    - SW-GRTs for Disease Prevention Research (Monica Taljaard, July 11, 2018)
    - Design and Analysis of IRGTs in Public Health (Sherri Pals, April 24, 2017)
    - Research Methods Resources for Clinical Trials Involving Groups or Clusters (David Murray, December 13, 2017)
- Research Methods Resources Website
  - <u>https://researchmethodsresources.nih.gov/</u>
  - Material on GRTs and IRGTs and a sample size calculator for GRTs.

![](_page_68_Picture_0.jpeg)

Health Care Systems Research Collaboratory

#### Demystifying Biostatistical Concepts for Embedded Pragmatic Clinical Trials

June 19, 2020

Elizabeth L. Turner, PhD, Duke University Patrick J. Heagerty, PhD, University of Washington David M. Murray, PhD, National Institutes of Health

Thank you

Any questions or comments?