

SGC_3A: The Insight Hypothesis

Amy Rae Fox

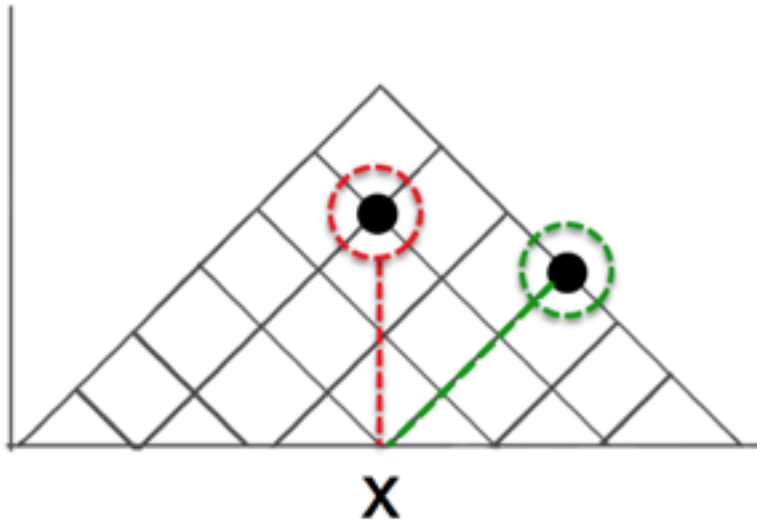
Contents

INTRODUCTION	1
Hypotheses	2
METHODS	2
Design	2
Sample	2
Materials	3
Procedure	3
DESCRIPTIVES	4
Participants	4
Response Accuracy	5
Response Latency	7
TODO ADD ITEM LEVEL	9
HYPOTHESIS TESTING	9
Response Accuracy by Condition	9
[EXPLORE]	9
[MODEL]	11
[REPLICATION]	12
DILLIGENCE	13
Assumptions of Wilcoxon Rank-Sum	13
DATA DICTIONARY	13
WIP PUBLIC WEBSITE VERSION OF ANALYSIS	

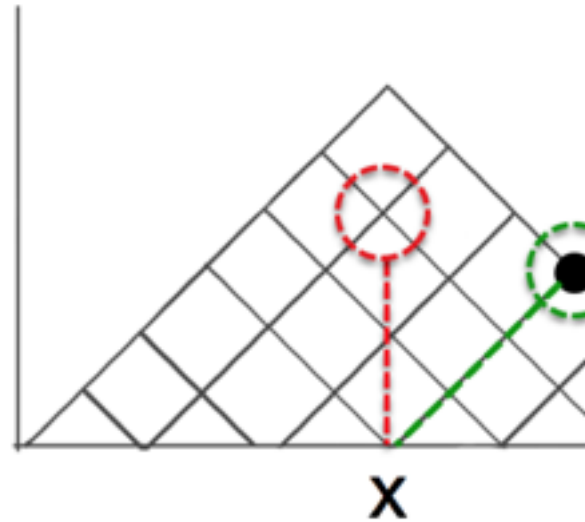
INTRODUCTION

In Study 3A we explore a hypothesis that emerged from analysis of Study 2, namely that presenting a learning with a situation that induces a state of impasse will increase the probability they have a moment of insight. In the context of Study 2, an impasse state was (unintentionally) induced when the combination of question + data set yielded no available answer in the incorrect (cartesian) interpretation of the graph. In Study 3A, we test this hypothesis by comparing performance between a (treatment) group receiving impasse-inducing questions followed by normal questions, and a non-impasse control.

control



impasse



Hypotheses

H1. Learners posed with impasse-inducing questions will be more likely to correct interpret the graph.

```
#FOR PUBLIC WEB VERSION
# ---
# **To try the study yourself: **
# visit TODO INSERT LINK
# *Enter "github" as your session code, and number of the condition you wish to test*
# session code= GITHUB
# condition code for CONTROL = 111
# condition code for IMPASSE = 121
# <br> <br>
```

METHODS

Design

We employed a mixed design with 1 between-subjects factor with 2 levels (Scaffold: control, impasse) and 15 items (within-subjects factor).

Independent Variables: B-S (Scaffold: control,impasse) W-S (Item x 15)

Dependent Variables 1. Response Accuracy : Is the response triangular-correct? 0 (false), 1 (true) 2. Response Latency : Time from stimulus onset to clicking 'Submit' button: time in (s)

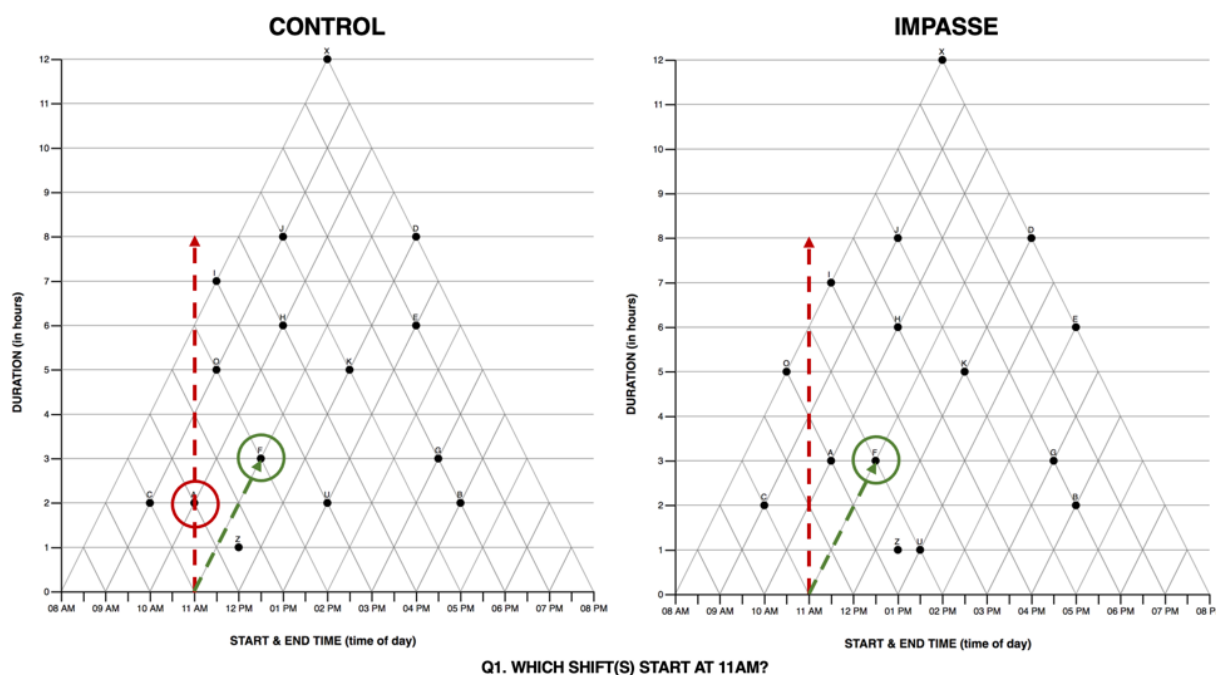
Sample

Data was collected by convenience sample of a university subject pool. Initial data (Fall 2017, Spring 2018) were collected in-person, with large groups of students simultaneously completing the study (independently)

in a computer lab. In Fall 2021 we collected additional data to replicate results in a remote format (students completing the study asynchronously on their own computers).

Materials

Stimuli consisted of a series of 15 graph comprehension questions, each testing a different combination of time interval relations, to be read from a Triangular-Model graph. In both conditions, the questions were identical. Participants in the IMPASSE condition saw a graph with a slightly different dataset, such that the some of the questions posed an IMPASSE structure, where no datapoint intersects the orthogonal projection from the x-axis required if reading the graph in accordance with the cartesian coordinate system. The complete list of questions can be found here. Examples of graphs+datasets stimuli for each condition are depicted below. The green line indicates the ideal-scanpath to the correct (triangular) answer to the first question, and the red line indicates the (incorrect) orthogonal interpretation. In the IMPASSE figure (at right), there are no data points that intersect the red line.



Procedure

Participants completed the study via a web-browser. Upon starting, they submitted informed consent, before reading task instructions. Participants were introduced to a scenario in which they were to play the role of a project manager, scheduling shifts for a group of employees. The schedule of the employees would be presented in a graph, and they would be answering question about the schedule. Then participants completed a test block of 15 items. In the IMPASSE condition, the first five questions included an IMPASSE problem state. The remaining 10 questions were not structured as impasse. Following the test block, participants answered a free-response question about their strategy for reading the graph, followed by a demographic questionnaire and debrief.

DESCRIPTIVES

```
#IMPORT PARTICIPANT DATA from fall and spring files
fall_participants <- "data/fall17_sgc3a_participants.csv"
spring_participants <- "data/spring18_sgc3a_participants.csv"
online_participants <- "data/fall21_sgc3a_participants.csv"

df_fall <- read.csv(fall_participants)
df_spring <- read.csv(spring_participants)
df_online <- read.csv(online_participants)

#indicate study modality
df_fall$mode <- "lab"
df_spring$mode <- "lab"
df_online$mode <- "online"

#Create combined data frame
df_subjects <- rbind(df_fall, df_spring, df_online) #, df_replication)
df_subjects$tri_min <- df_subjects$triangular_time / 1000 / 60
df_subjects$test_min <- df_subjects$tt_t / 1000 / 60
df_subjects$learn_min <- df_subjects$ts_t / 1000 / 60

#Create factors
df_subjects <- df_subjects %>% mutate(
  subject = as.factor(subject),
  session = as.factor(session),
  term = as.factor(term),
  condition = as.factor(condition),
  explicit = as.factor(explicit),
  impasse = as.factor(impasse),
  axis = as.factor(axis)
)

df_online <- df_subjects %>% filter(term == "fall21")
df_lab <- df_subjects %>% filter(term != "fall21")

#Remove extraneous dfs
rm(df_fall, df_spring)
```

Participants

```
#Describe participants
subject.stats <- rbind(
  "lab" = df_lab %>% dplyr::select(age) %>% unlist() %>% favstats(),
  "online" = df_online %>% dplyr::select(age) %>% unlist() %>% favstats()
)

subject.stats$female <- c(
  (df_lab %>% filter(sex=="Female") %>% count())$n,
  (df_online %>% filter(sex=="Female") %>% count())$n
)

#participants per condition
```

```
mode.stats <- df_subjects %>% group_by(mode, condition) %>%
  summarize(n=n())
```

For **in-person** collection, 126 participants (60 % female) undergraduate STEM majors at a public American University participated *in person* in exchange for course credit (age: 18 - 33 years). Participants were randomly assigned to one of two experimental groups, with 62 in the control condition, and 64 in the experimental IMPASSE condition.

For **online replication** 140 participants (70 % female) undergraduate STEM majors at a public American University participated *online, asynchronously* in exchange for course credit (age: 18 - 31 years). Participants were randomly assigned to one of two experimental groups, with 68 in the control condition, and 72 in the experimental IMPASSE condition.

Response Accuracy

Response accuracy refers to how many questions the subject answers with a correct (triangular) interpretation.

```
#DESCRIBE distribution of triangular-correct scores
score.stats <- rbind(
  "lab"= favstats(df_lab$triangular_score),
  "online"= favstats(df_online$triangular_score)
)
score.stats
```

##		min	Q1	median	Q3	max	mean	sd	n	missing
##	lab	1	2	3	11	15	5.809524	4.893611	126	0
##	online	1	2	2	8	15	5.007143	4.738479	140	0

For *in person* collection, accuracy scores (n = 126) range from 1 to 15 with a mean score of (M = 5.81, SD = 4.89).

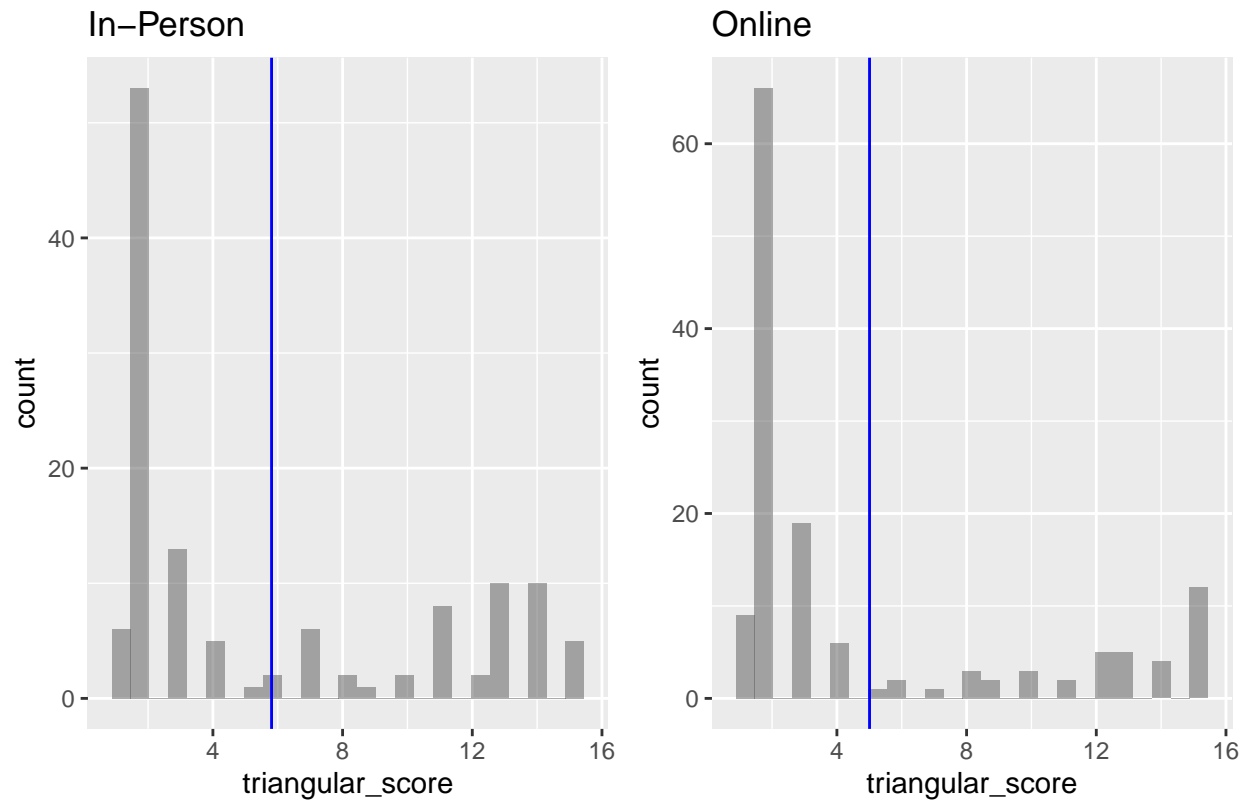
For *online replication*, (online) accuracy scores (n = 140) range from 1 to 15 with a mean score of (M = 5.01, SD = 4.74).

```
#VISUALIZE distribution of response accuracy
plab <- gf_histogram(~triangular_score, data = df_lab) %>%
  gf_vline(xintercept = score.stats["lab",]$mean, color = "blue") +
  labs(title="In-Person")

ponline <- gf_histogram(~triangular_score, data = df_online) %>%
  gf_vline(xintercept = score.stats["online",]$mean, color = "blue") +
  labs(title="Online")

plot <- ggarrange(plab, ponline, common.legend = TRUE, nrow = 1, ncol = 2)
annotate_figure(plot, top = text_grob("Score Accuracy by Study",
  color = "black", face = "bold", size = 14))
```

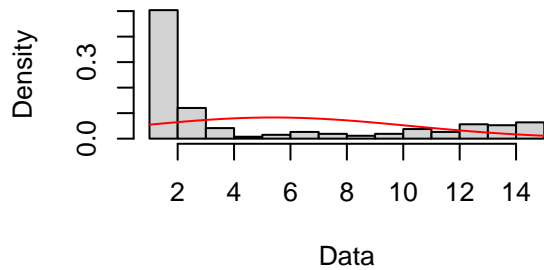
Score Accuracy by Study



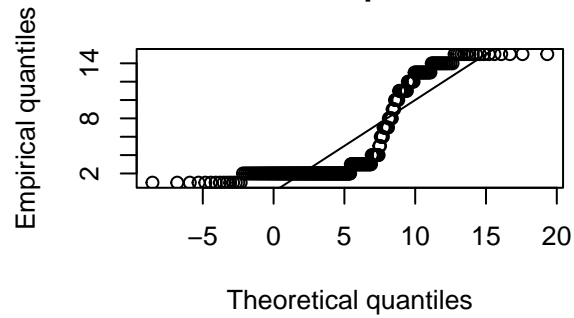
However, inspection of the quantile-quantile plots reveal that response accuracy does not approximate a normal distribution. After exploring several transformations and comparing against alternative distributions (log-normal, poisson, exp, nbinom), we conclude that we will need to use robust tests to analyze response accuracy.

```
plot(fitdist(df_subjects$triangular_score, "norm"))
```

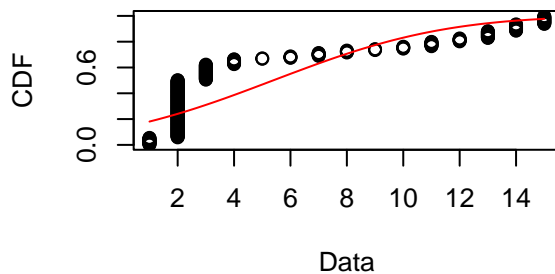
Empirical and theoretical dens.



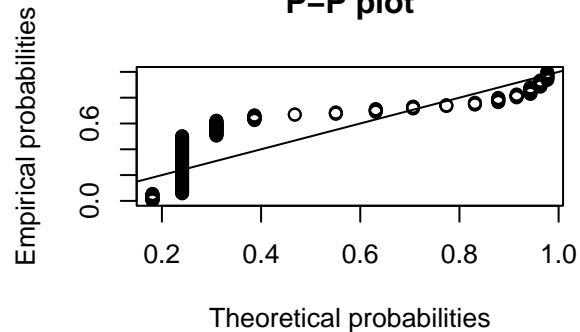
Q-Q plot



Empirical and theoretical CDFs



P-P plot



```
#EXPLORE ALTERNATIVES
# fit_n <- fitdist(df_fall$triangular_score, "norm")
# fit_p <- fitdist(df_fall$triangular_score, "pois")
# fit_b <- fitdist(df_fall$triangular_score, "nbinom")

# par(mfrow=c(2,2))
# plot.legend <- c("normal", "poisson", "nbinomial")
# denscomp(list(fit_n, fit_p, fit_b), legendtext = plot.legend)
# cdfcomp(list(fit_n, fit_p, fit_b), legendtext = plot.legend)
# qqcomp(list(fit_n, fit_p, fit_b), legendtext = plot.legend)
# ppcomp(list(fit_n, fit_p, fit_b), legendtext = plot.legend)
```

Response Latency

```
#DESCRIBE distribution of response time
time.stats <- rbind(
  "lab"= favstats(df_lab$tri_min),
  "online"= favstats(df_online$tri_min)
)
time.stats <- time.stats %>% dplyr::select(-missing) #don't need missing column
time.stats
```

	min	Q1	median	Q3	max	mean	sd	n
lab	3.770800	7.283617	8.804333	10.71089	19.98955	9.253020	2.902315	126
online	1.983417	7.083704	8.811267	11.84312	27.87285	9.924459	4.527250	140

For *in person* response latency (for test block) (n = 126) range from 3.77 to 19.99 minutes, with a mean duration of (M = 9.25, SD = 2.9) minutes.

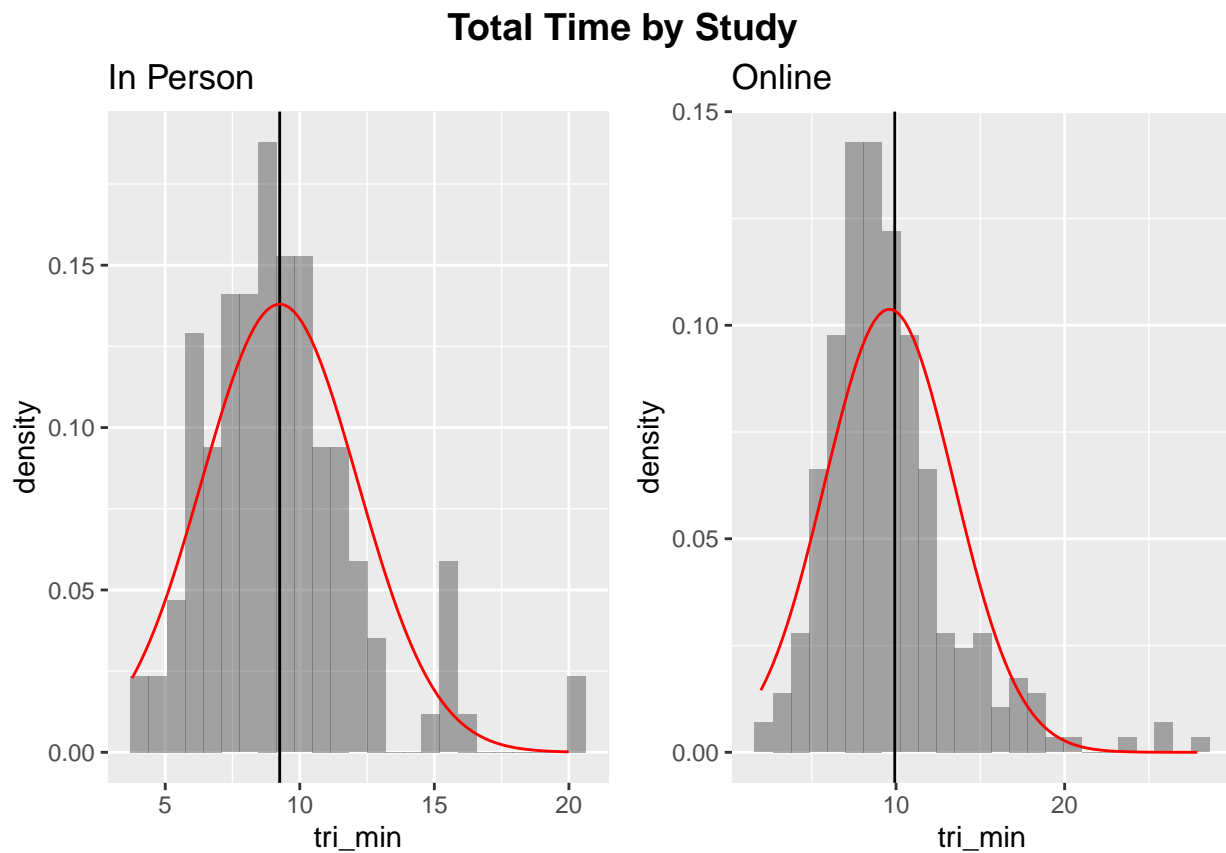
For *online replication* (online) response latency (for test block) ($n = 140$) range from 1.98 to 27.87 minutes, with a mean duration of ($M = 9.92$, $SD = 4.53$).

```
#VISUALIZE distribution of response time
plab <- gf_dhistogram(~tri_min, data = df_lab) %>%
  gf_vline(xintercept = time.stats["lab",]$mean, color = "black") %>%
  gf_fitdistr(color="red")+
  labs(title="In Person")

ponline <- gf_dhistogram(~tri_min, data = df_subjects) %>%
  gf_vline(xintercept = time.stats["online",]$mean, color = "black") %>%
  gf_fitdistr(color="red")+
  labs(title="Online")

plot <- ggarrange(plab, ponline, common.legend = TRUE, nrow = 1, ncol = 2)

annotate_figure(plot, top = text_grob("Total Time by Study",
  color = "black", face = "bold", size = 14))
```



The data may need to be log-transformed. But we will address this when modelling with the variable.

TODO ADD ITEM LEVEL

HYPOTHESIS TESTING

Response Accuracy by Condition

The experimental hypothesis (H1) is that structuring the data to pose an impasse (condition 121) will produce significantly better performance than non-impasse (condition 111). The null hypothesis (H0) is that there will be no difference in performance between conditions.

[EXPLORE]

```
#DESCRIBE scores by condition
score.cond.stats <- rbind(
  "lab" = favstats(triangular_score ~ condition, data = df_lab),
  "online" = favstats(triangular_score ~ condition, data = df_online)
)
score.cond.stats
```

##	condition	min	Q1	median	Q3	max	mean	sd	n	missing
## lab.1	111	1	2	2.0	3.75	15	4.500000	4.643875	62	0
## lab.2	121	1	2	6.5	12.00	15	7.078125	4.828174	64	0
## online.1	111	1	2	2.0	3.00	15	3.897059	3.996789	68	0
## online.2	121	1	2	3.0	11.25	15	6.055556	5.156396	72	0

For **in person** study, participants in the impasse group had (on average) higher scores ($M = 7.08$ $SD = 4.83$) than those in the non-impasse control group ($M = 4.5$, $SD = 4.64$).

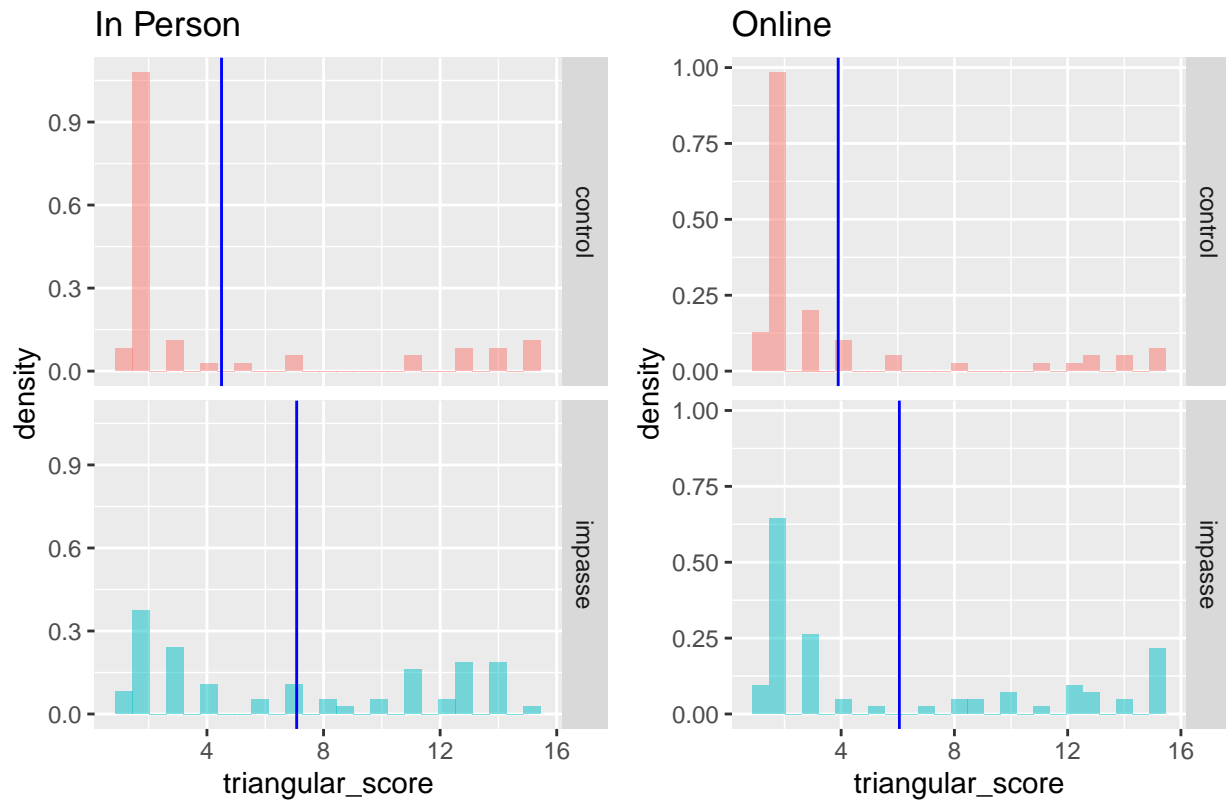
For **online replication** study, participants in the impasse group had (on average) higher scores ($M = 6.06$ $SD = 5.16$) than those in the non-impasse control group ($M = 3.9$, $SD = 4$).

```
#VISUALIZE scores by condition
condlables <- c("111"="control", "121"="impasse")
plab <- gf_dhistogram( ~triangular_score, fill= ~condition, data = df_lab) %>%
  gf_facet_grid(condition~., labeller=labeller(condition=condlables)) %>%
  gf_vline(xintercept = ~mean, data = score.cond.stats[c(1:2),], color = "blue")+
  labs(title="In Person")

ponline <- gf_dhistogram( ~triangular_score, fill= ~condition, data = df_online) %>%
  gf_facet_grid(condition~., labeller=labeller(condition=condlables)) %>%
  gf_vline(xintercept = ~mean, data = score.cond.stats[c(3:4),], color = "blue")+
  labs(title="Online")

plot <- ggarrange(plab, ponline, legend = FALSE, nrow = 1, ncol =2)
annotate_figure(plot, top = text_grob("Score Accuracy by Condition",
  color = "black", face = "bold", size = 14))
```

Score Accuracy by Condition

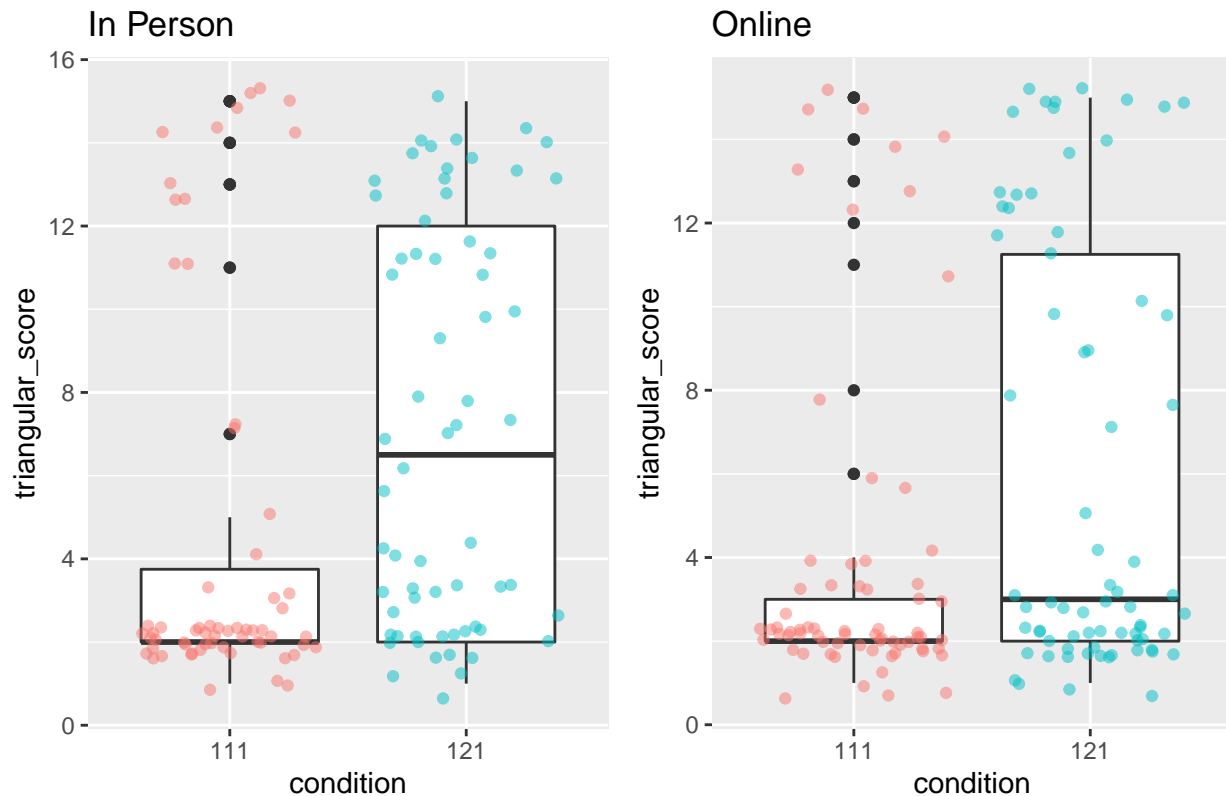


```
#VISUALIZE scores by condition
plab <- gf_boxplot(triangular_score ~ condition, data=df_lab) %>%
  gf_jitter(color=~condition, alpha=0.5) +
  labs (title = "In Person")

ponline <- gf_boxplot(triangular_score ~ condition, data = df_online) %>%
  gf_jitter(color=~condition, alpha=0.5)+
  labs(title ="Online")

plot <-ggarrange(plab, ponline, legend = FALSE, nrow = 1, ncol =2)
annotate_figure(plot, top = text_grob("Score Accuracy by Condition",
  color = "black", face = "bold", size = 14))
```

Score Accuracy by Condition



[MODEL]

Because the response accuracy data are not-normal, we will test the veracity of our hypothesis using Wilcoxon rank-sum test (Wilcoxon, 1945).

```
#Wilcoxon Rank Sum Test
m1 <- wilcox.test(triangular_score ~ condition, data = df_lab, exact=FALSE,
                  alternative="less")
m1
```

```
##
## Wilcoxon rank sum test with continuity correction
##
## data: triangular_score by condition
## W = 1284.5, p-value = 0.0001922
## alternative hypothesis: true location shift is less than 0
```

```
#Calculate Effect size
m1eff <- wilcox_effsize(triangular_score ~ condition, data = df_lab)
m1eff
```

```
## # A tibble: 1 x 7
##   .y.      group1 group2 effsize    n1    n2 magnitude
## * <chr>    <chr> <chr>   <dbl> <int> <int> <ord>
## 1 triangular_score 111   121    0.317   62   64 moderate
```

In Person A one-tailed Wilcoxon rank-sum test reveals that the cumulative score (number of triangular-consistent responses)(Mdn = 6.5 points) in the IMPASSE group were significantly higher than than scores

in the non-impasse control condition (Mdn = 2), $W = 1284.5$, $p < 0.001$, $r = 0.3165359$, a moderate-sized effect.

```
#Wilcoxon Rank Sum Test
m1 <- wilcox.test(triangular_score ~ condition, data = df_online, exact=FALSE,
                  alternative="less")
m1

##
## Wilcoxon rank sum test with continuity correction
##
## data: triangular_score by condition
## W = 1848.5, p-value = 0.004087
## alternative hypothesis: true location shift is less than 0

#Calculate Effect size
m1eff <- wilcox_effsize(triangular_score ~ condition, data = df_online)
m1eff

## # A tibble: 1 x 7
##   .y.      group1 group2 effsize    n1    n2 magnitude
## * <chr>      <chr> <chr>   <dbl> <int> <int> <ord>
## 1 triangular_score 111    121    0.224    68    72 small
```

Online A one-tailed Wilcoxon rank-sum test reveals that the cumulative score (number of triangular-consistent responses)(Mdn = 3 points) in the IMPASSE group were significantly higher than than scores in the non-impasse control condition (Mdn = 2), $W = 1848.5$, $p < 0.05$, $r = 0.2237108$, a small effect.

TODO:: Investigate outliers in online study; may need to have stricter response time + strategy-consistent response criteria

[REPLICATION]

Because the effect of impasse scaffold was smaller in the remote-online sample than in-person sample, we'll test whether scores in in-person sample were significantly smaller than those online using another Wilcoxon rank-sum test.

```
temp <- df_subjects %>% filter(condition==121)

#Wilcoxon Rank Sum Test
m1 <- wilcox.test(triangular_score ~ mode, data = temp, exact=FALSE,
                  alternative="greater")
m1

##
## Wilcoxon rank sum test with continuity correction
##
## data: triangular_score by mode
## W = 2620.5, p-value = 0.08069
## alternative hypothesis: true location shift is greater than 0

#Calculate Effect size
m1eff <- wilcox_effsize(triangular_score ~ mode, data = temp)
m1eff

## # A tibble: 1 x 7
##   .y.      group1 group2 effsize    n1    n2 magnitude
## * <chr>      <chr> <chr>   <dbl> <int> <int> <ord>
```

```
## 1 triangular_score lab    online    0.120    64    72 small
```

A one-tailed Wilcoxon rank-sum test reveals that the cumulative score (number of triangular-consistent responses)(Mdn = 6.5 points) in the IN-PERSON (impasse condition) sample *were not* significantly higher than than scores in the REMOTE-ONLINE (impasse condition) sample (Mdn = 3), $W = 2620.5$, $p = 0.08$ $r = 0.1202775$.

Inference Our replication study had comparable results to the in-person study, however the effect size was smaller, therefore we should consider increasing the sample size of future studies conducted online so as to ensure we have sufficient statistical power detect a smaller effect.

DILLIGENCE

Assumptions of Wilcoxon Rank-Sum

The Wilcoxon rank-sum test is the non-parametric alternative to a independent samples t-test and requires the following assumptions:

1. The two samples are independent of one another <- MET by random sampling + assignment
2. The two populations have equal variance or spread <-TEST

```
df_subjects %>% group_by(condition,mode) %>% summarize(var=var(triangular_score))
```

```
## # A tibble: 4 x 3
## # Groups:   condition [2]
##   condition mode    var
##   <fct>      <chr> <dbl>
## 1 111      lab     21.6
## 2 111    online    16.0
## 3 121      lab     23.3
## 4 121    online    26.6
```

Variances between conditions in both Lab and Online samples are comparable

DATA DICTIONARY