

# Data types and encodings

# Data types

We can classify data types into two broad categories

- Ordered
- Categorical

# Ordered data types

All ordered data has an *implicit* ordering

- Ordinal vs. quantitative
- Sequential and diverging
- Temporal

# Ordinal data types

There is a meaningful ordering, but we cannot do full-fledged arithmetic

- Shirt sizes (XS, S, M, L, XL)
- Rankings (1st, 2nd, 3rd, 4th, 5th)

# Quantitative data

There is a meaningful ordering, and arithmetic operations are supported

- Lengths, heights, etc...
- Counts
- Prices
- ....

# Sequential and diverging data

*sequential*: where there is a homogeneous range from a minimum to a maximum value

*diverging*: can be deconstructed into two sequences pointing in opposite directions that meet at a *common zero* point

- Sequential: height of mountains, depth of seas
- Diverging: the elevation of points on earth surface; the PH scale (acid to alkaline, with middle point on 7.5)

---

This distinction will be useful mainly for defining meaningful scales

# Temporal data

Time is a complex subject, so temporal data requires special care

- Does every year have 365 days?
- Does every day have 24 hours?
- Does every minute have 60 seconds?
- Leap years
- Daylight Saving Time
- Leap seconds

Dates and times are hard because they have to reconcile two physical phenomena (the rotation of the Earth and its orbit around the sun) with a whole raft of geopolitical phenomena including months, time zones, and DST.

# Lubridate

The `lubridate` package provides a series of facilities to deal with date and times.

It takes care for you of all the gory details of time handling on computers, including time zones, leap years/seconds, daylight saving times, etc...

```
1 renv::install("lubridate")  
2 library(lubridate)
```

It provides three datatypes

- `date`
- `time`
- `datetime`



# Creating dates and datetimes

## Creating dates

```
1 ymd("2019-10-25")
```

```
[1] "2019-10-25"
```

```
1 mdy("October 25th, 2019")
```

```
[1] "2019-10-25"
```

```
1 dmy("25-Oct-2019")
```

```
[1] "2019-10-25"
```

```
1 # On dates that cannot be parsed will return NA  
2 ymd("2019-4r-31")
```

```
[1] NA
```

## Creating datetimes

```
1 ymd_hms("2019-10-25 08:10:03")
```

```
[1] "2019-10-25 08:10:03 UTC"
```

## From components

```
1 make_date(2019, 10, 25)
```

```
[1] "2019-10-25"
```

```
1 make_datetime(2019, 10, 25, 8, 10, 25)
```

```
[1] "2019-10-25 08:10:25 UTC"
```

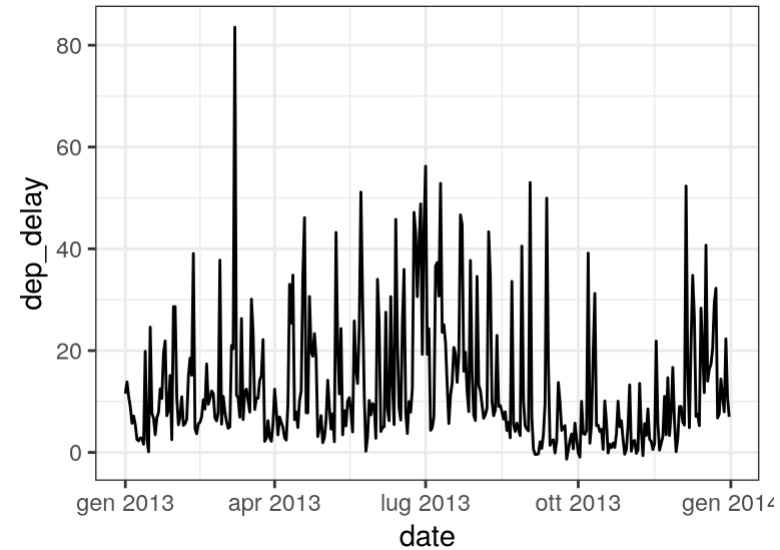
## With vectorization

```
1 dates <- c(  
2   "October 25th, 2019",  
3   "October 26th, 2019",  
4   "Novembber 3rd, 2019",  
5   "November 4rd, 2019"  
6 )  
7 mdy(dates)
```

```
[1] "2019-10-25" "2019-10-26" NA  
"2019-11-04"
```

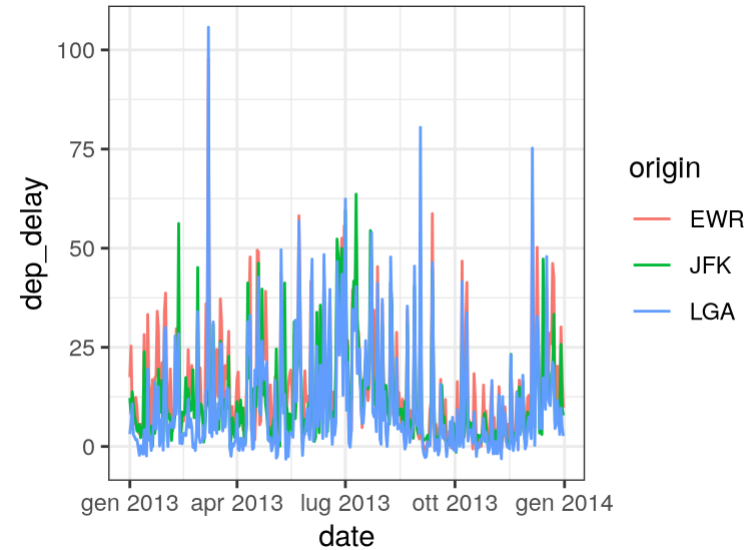
# Bringing proper dates to nycflights13

```
1 library(nycflights13)
2 library(tidyverse)
3 library(lubridate)
4
5 flights |>
6   drop_na(dep_delay) |>
7   mutate(date = make_date(year, month, day)) |>
8   group_by(date) |>
9   summarise(dep_delay = mean(dep_delay)) |>
10  ggplot(aes(x=date, y=dep_delay)) +
11    geom_line()
```



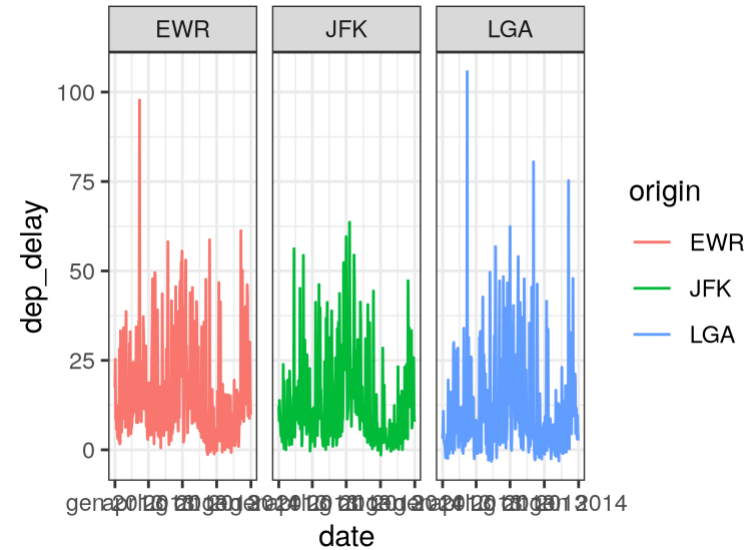
# Bringing proper dates to nycflights13

```
1 library(nycflights13)
2 library(tidyverse)
3 library(lubridate)
4
5 flights |>
6   drop_na(dep_delay) |>
7   mutate(date = make_date(year, month, day)) |>
8   group_by(date, origin) |>
9   summarise(dep_delay = mean(dep_delay)) |>
10  ggplot(aes(x=date, y=dep_delay, color=origin)) +
11    geom_line()
```



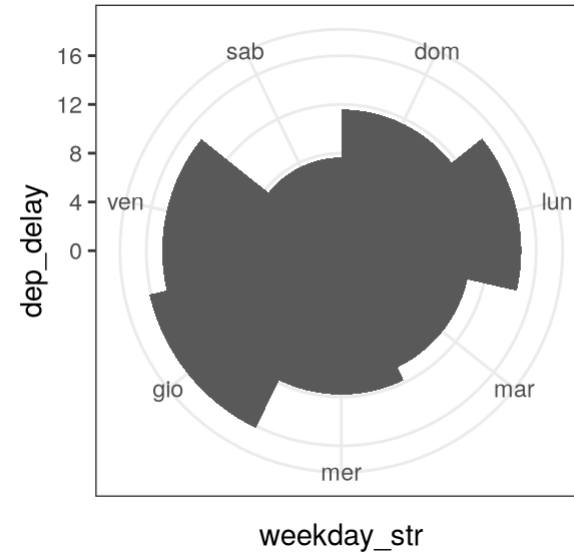
# Bringing proper dates to nycflights13

```
1 library(nycflights13)
2 library(tidyverse)
3 library(lubridate)
4
5 flights |>
6   drop_na(dep_delay) |>
7   mutate(date = make_date(year, month, day)) |>
8   group_by(date, origin) |>
9   summarise(dep_delay = mean(dep_delay)) |>
10  ggplot(aes(x=date, y=dep_delay, color=origin)) +
11    geom_line() +
12    facet_wrap(vars(origin))
```



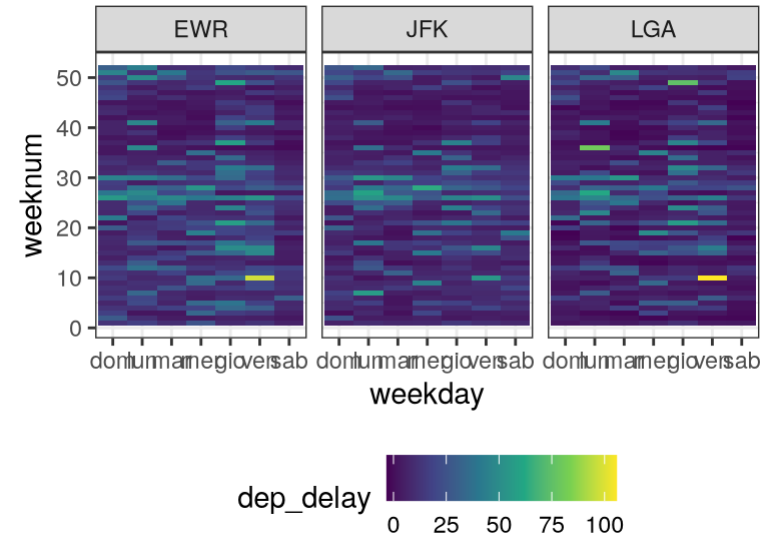
# Bringing proper dates to nycflights13

```
1 flights |>
2   drop_na(dep_delay) |>
3   mutate(
4     date = make_date(year, month, day),
5     weekday = wday(date, label=FALSE), # Gets the d
6     weekday_str = wday(date, label=TRUE)
7   ) |>
8   group_by(weekday, weekday_str) |>
9   summarise(dep_delay = mean(dep_delay)) |>
10  ggplot(aes(x=weekday_str, y=dep_delay)) +
11    geom_col(width = 1) +
12    coord_polar()
```



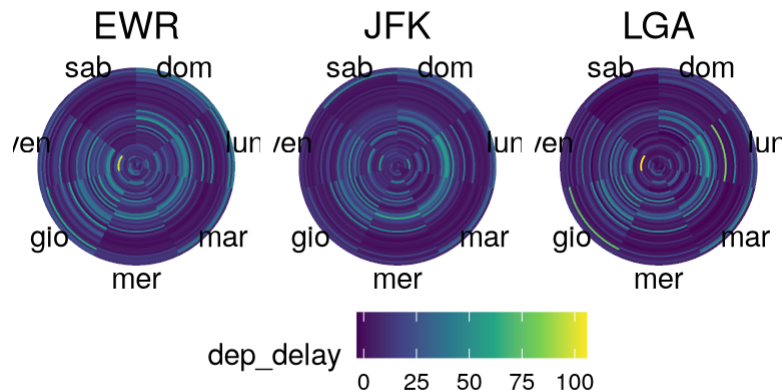
# Bringing proper dates to nycflights13

```
1 flights |>
2   drop_na(dep_delay) |>
3   mutate(date = make_date(year, month, day)) |>
4   group_by(date, origin) |>
5   summarise(dep_delay = mean(dep_delay)) |>
6   mutate(
7     weekday = wday(date, label=TRUE),
8     weeknum = isoweek(date)) |>
9   ggplot(aes(x=weekday, y=weeknum,
10             fill=dep_delay)) +
11   geom_tile() +
12   scale_fill_continuous(type="viridis") +
13   facet_wrap(vars(origin)) +
14   theme_bw() +
15   theme(legend.position = 'bottom')
```



# Bringing proper dates to nycflights13

```
1 flights |>
2   drop_na(dep_delay) |>
3   mutate(date = make_date(year, month, day)) |>
4   group_by(date, origin) |>
5   summarise(dep_delay = mean(dep_delay)) |>
6   mutate(
7     weekday = wday(date, label=TRUE),
8     weeknum = isoweek(date)) |>
9   ggplot(aes(x=weekday, y=weeknum,
10             fill=dep_delay, color=dep_delay)) +
11   geom_tile() +
12   scale_fill_continuous(type="viridis", aesthetics=
13   facet_wrap(vars(origin)) +
14   coord_polar() +
15   theme_void() +
16   theme(axis.text.x = element_text(),
17         strip.text = element_text(size=14),
18         legend.position = 'bottom')
```



# Categorical

Categorical types do not have an implicit ordering.

With categories you can just compare for equality

- Apples, oranges, bananas, grapes...

You can enforce an explicit ordering when needed



# Forcats

This is a tidyverse library that helps with working with categorical data, which in R is called a `factor`

```
1 library(forcats)
```

# Creating factors

You can use the `factor` function to create factors, which is vectorized

```
1 library(gapminder)
2
3 mutate(gapminder,
4         country = factor(country))
```

```
# A tibble: 1,704 × 6
  country      continent  year lifeExp      pop gdpPercap
  <fct>        <fct>    <int> <dbl>    <int> <dbl>
1 Afghanistan Asia      1952  28.8  8425333  779.
2 Afghanistan Asia      1957  30.3  9240934  821.
3 Afghanistan Asia      1962  32.0 10267083  853.
4 Afghanistan Asia      1967  34.0 11537966  836.
5 Afghanistan Asia      1972  36.1 13079460  740.
6 Afghanistan Asia      1977  38.4 14880372  786.
7 Afghanistan Asia      1982  39.9 12881816  978.
8 Afghanistan Asia      1987  40.8 13867957  852.
9 Afghanistan Asia      1992  41.7 16317921  649.
10 Afghanistan Asia      1997  41.8 22227415  635.
# i 1,694 more rows
```

# Factors and strings

What is the difference between factors and strings?

- Factors can take only a pre-specified number of values

```
1 x <- factor(c("a", "b", "c"))
2 x[4] <- "d"
3 x
```

```
[1] a    b    c    <NA>
Levels: a b c
```

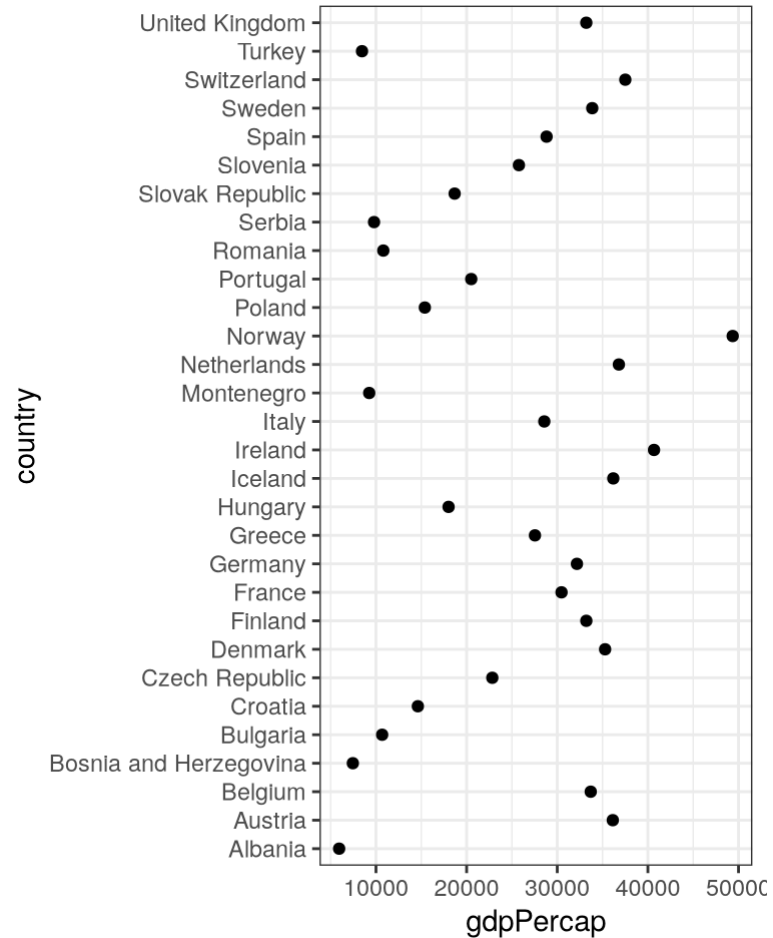
```
1 x <- factor(c("a", "b", "c"), levels = c("a", "b", "c", "d"))
2 x[4] <- "d"
3 x
```

```
[1] a b c d
Levels: a b c d
```

- You can impose an ordering different than alphabetical to factors.

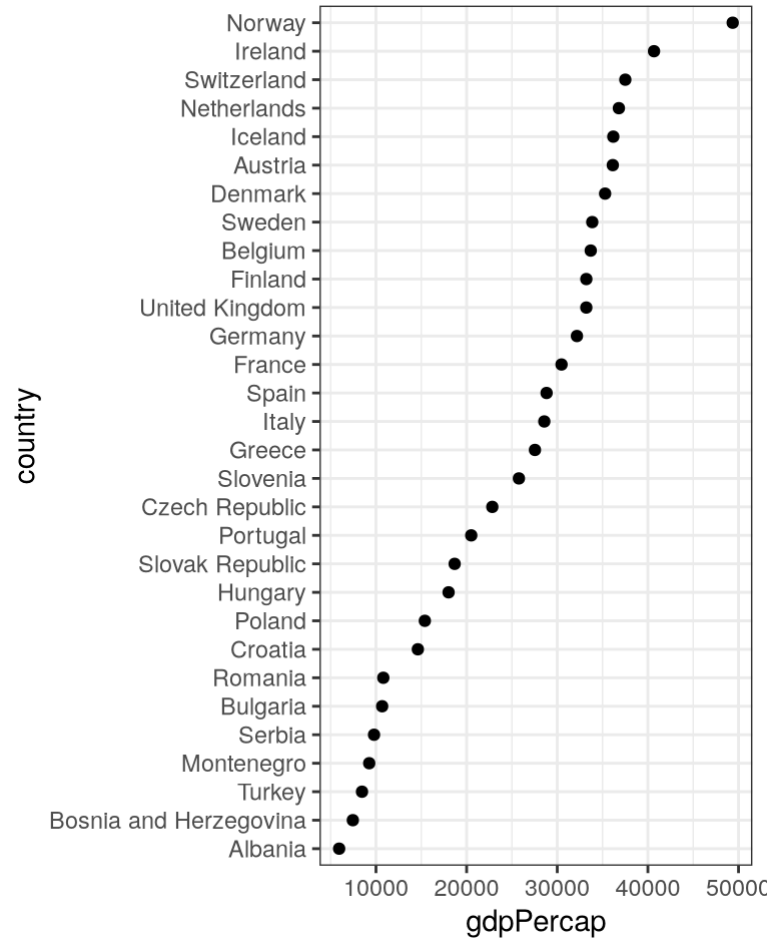
# Reordering factors

```
1 gapminder |>  
2   filter(year==2007) |>  
3   filter(continent == 'Europe') |>  
4   ggplot(aes(x=country, y=gdpPercap)) +  
5     geom_point() +  
6     coord_flip()
```



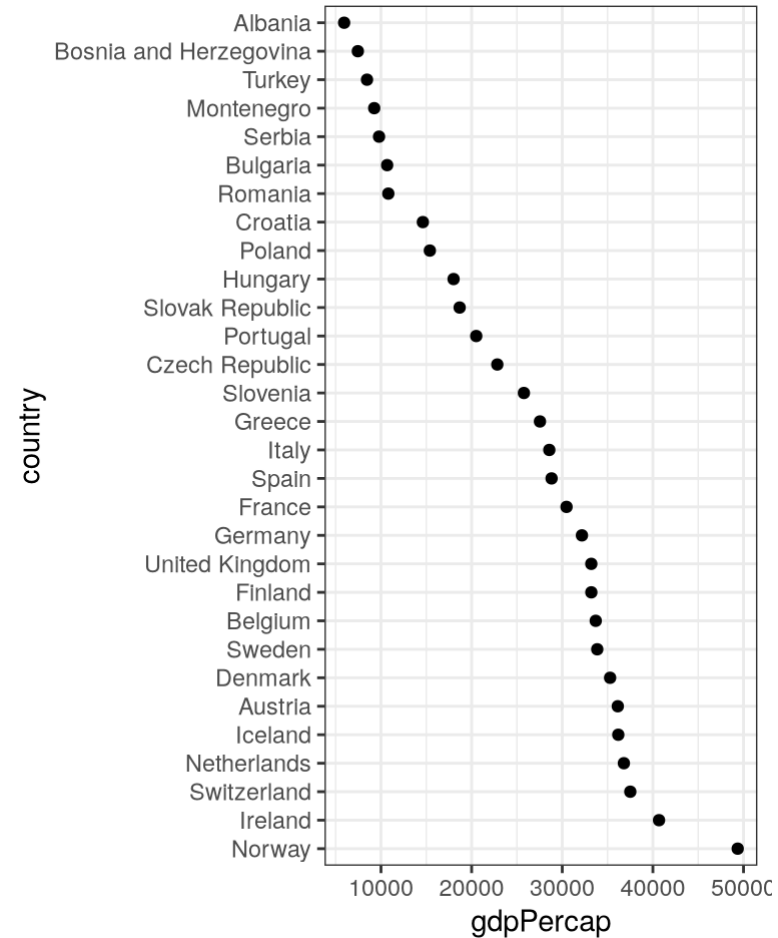
# Reordering factors

```
1 gapminder |>
2   filter(year==2007) |>
3   filter(continent == 'Europe') |>
4   mutate(
5     country = fct_reorder(country, gdpPercap)
6   ) |>
7   ggplot(aes(x=country, y=gdpPercap)) +
8     geom_point() +
9     coord_flip()
```



# Reordering factors

```
1 gapminder |>
2   filter(year==2007) |>
3   filter(continent == 'Europe') |>
4   mutate(
5     country = fct_reorder(country, desc(gdpPercap))
6   ) |>
7   ggplot(aes(x=country, y=gdpPercap)) +
8     geom_point() +
9     coord_flip()
```



# Aesthetics, geoms, and data

Care must be used when selecting the mapping between attributes and aesthetics, as well as when choosing the geometric objects we use.

# Aesthetic types

There are mainly two types of aesthetics

- Magnitude aesthetics
- Identity aesthetics



# Aesthetic types

There are mainly two types of aesthetics

- Magnitude aesthetics: ordered data
- Identity aesthetics: categorical data

➔ **Magnitude Channels: Ordered Attributes**

Position on common scale 

Position on unaligned scale 

Length (1D size) 

Tilt/angle 

Area (2D size) 

Depth (3D position) 

Color luminance 

Color saturation 

Curvature 

Volume (3D size) 

Same

Same

## → Identity Channels: Categorical Attributes

Spatial region



Color hue



Motion



Shape



# Properties

- Expressiveness: what can be encoded by what
  - Effectiveness: the importance should match the salience of the aesthetic, i.e. its noticeability
- 

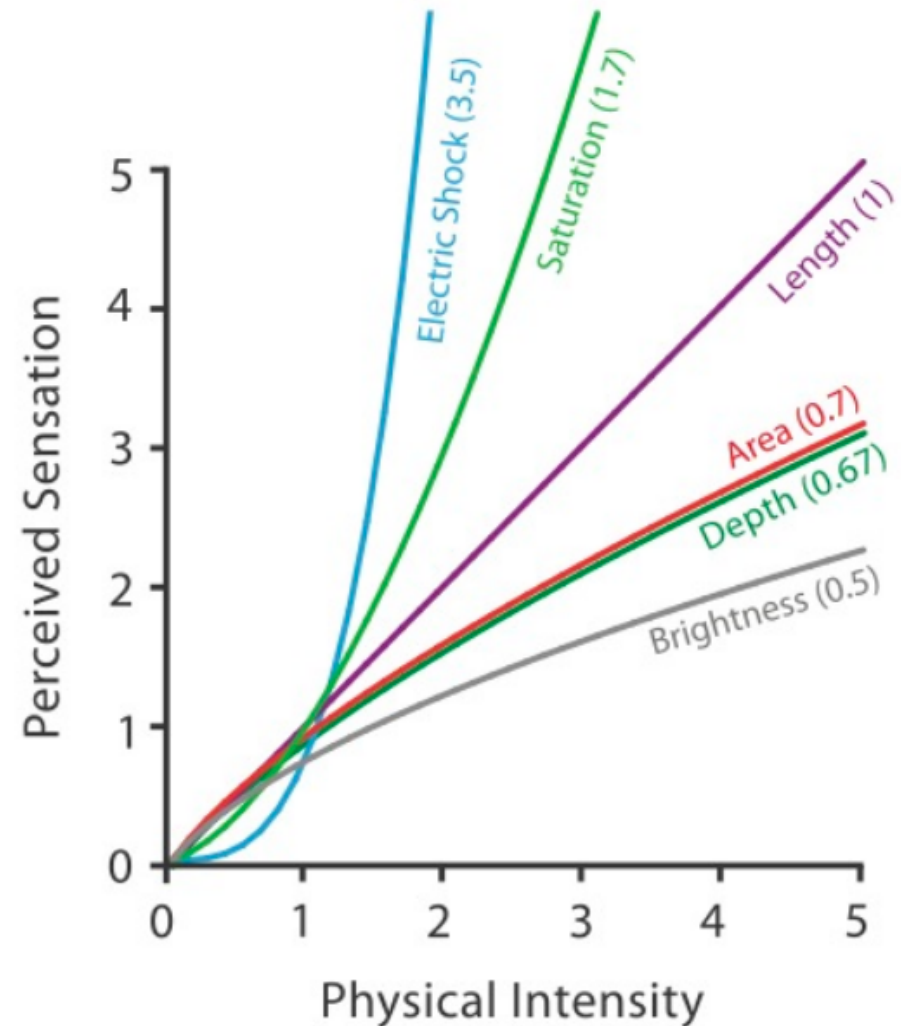
## When choosing a mapping

- Accuracy
- Discriminability
- Separability
- Popout

# Accuracy

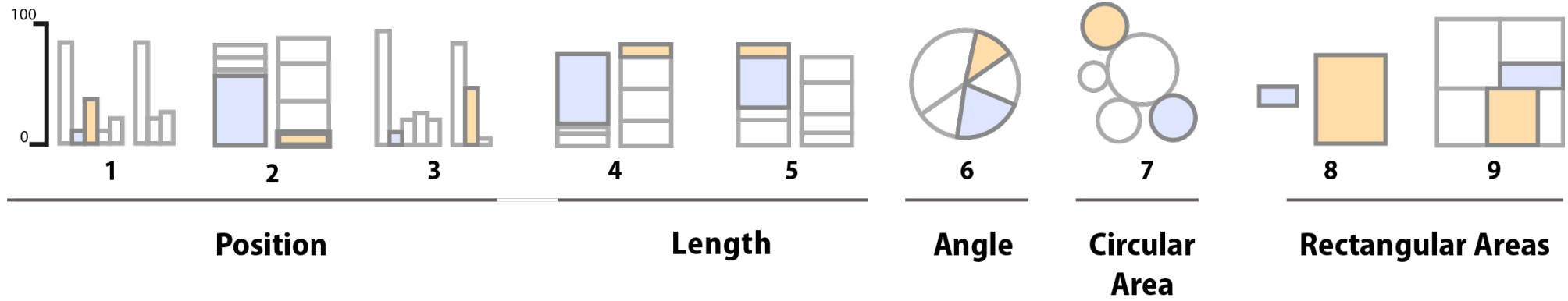
- Given a stimulus how close to the it measurement is the human perceptual judgement
- Stevens psychophysical power law [1975]
- Be careful with encoding linear information with a non perceptually-linear aesthetic

Steven's Psychophysical Power Law:  $S = I^N$

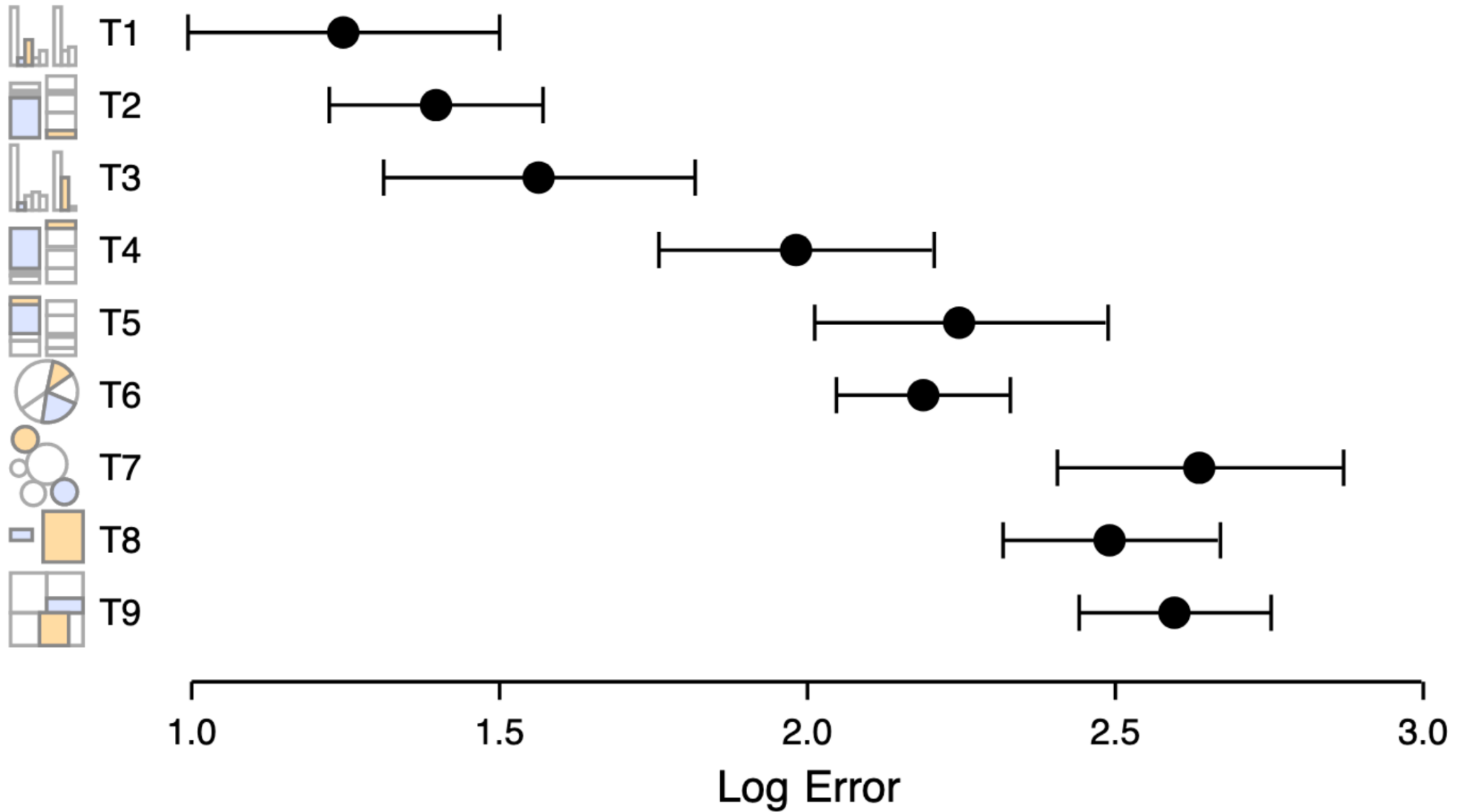




# Accuracy



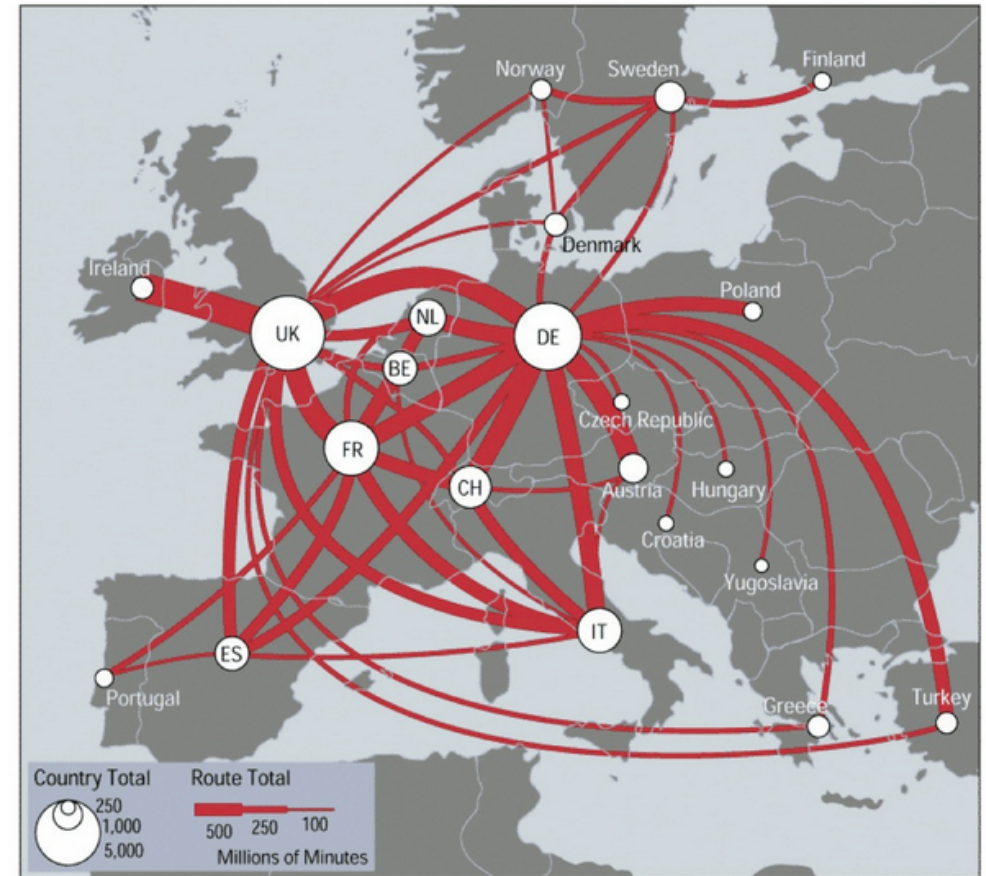
# Crowdsourced Results





# Discriminability

- If you encode data using a particular aesthetic, are the differences between items perceptible to the human as intended?
- How many levels can be distinguished?
- We should quantify the number of bins that are available in a visual channel



# Separability

Separable

Integral



color x location

color x shape

x-size x y-size

color x motion

size x orientation

r-g x y-b

# Popout

**Table 1. Our sales grew to \$600 million this year**

	Q1	Q2	Q3	Q4
Bob	26	35	72	84
Ellie	22	15	61	35
Gerrie	19	20	71	55
Jack	22	95	13	64
Jon	83	62	46	48
Karen	30	65	98	82
Ken	38	28	45	71
Lauren	98	81	41	63
Steve	16	50	23	41
Valerie	46	24	30	57
<b>Total</b>	<b>\$400</b>	<b>\$475</b>	<b>\$500</b>	<b>\$600</b>

# Popout

Table 1. Our sales grew to \$600 million this year

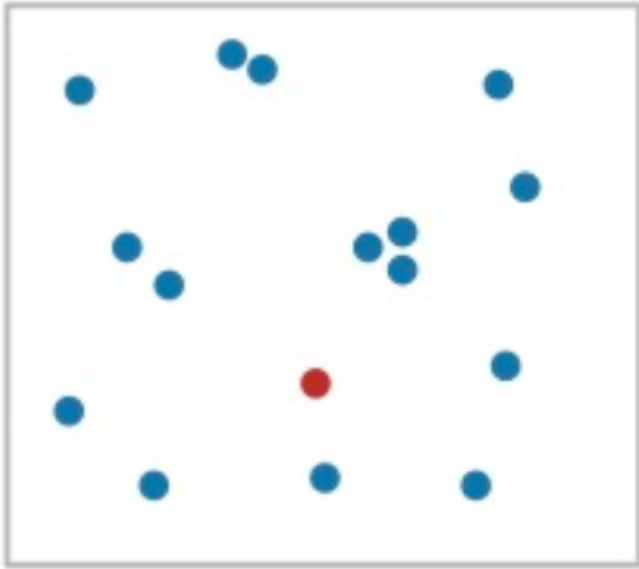
	Q1	Q2	Q3	Q4
Bob	26	35	72	84
Ellie	22	15	61	35
Gerrie	19	20	71	55
Jack	22	95	13	64
Jon	83	62	46	48
Karen	30	65	98	82
Ken	38	28	45	71
Lauren	98	81	41	63
Steve	16	50	23	41
Valerie	46	24	30	57
<b>Total</b>	<b>\$400</b>	<b>\$475</b>	<b>\$500</b>	<b>\$600</b>

# Popout

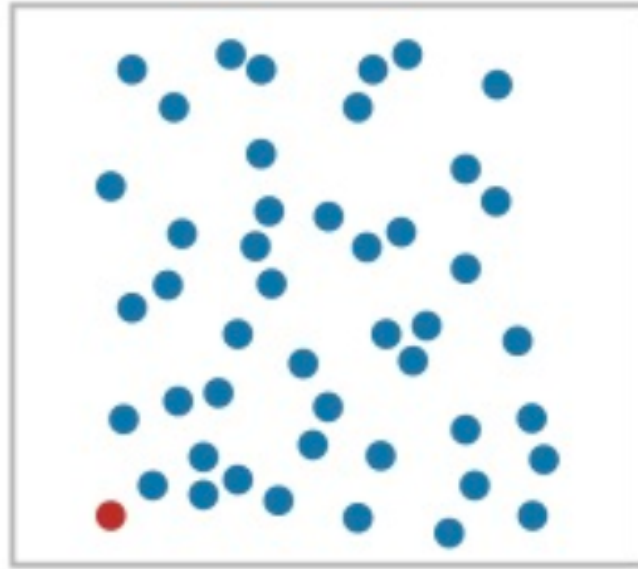
Table 1. Our sales grew to \$600 million this year

	Q1	Q2	Q3	Q4
Bob	26	35	72	<b>84</b>
Ellie	22	15	61	35
Gerrie	19	20	71	55
Jack	22	<b>95</b>	13	64
Jon	83	62	46	48
Karen	30	65	<b>98</b>	82
Ken	38	28	45	71
Lauren	<b>98</b>	81	41	63
Steve	16	50	23	41
Valerie	46	24	30	57
<b>Total</b>	<b>\$400</b>	<b>\$475</b>	<b>\$500</b>	<b>\$600</b>

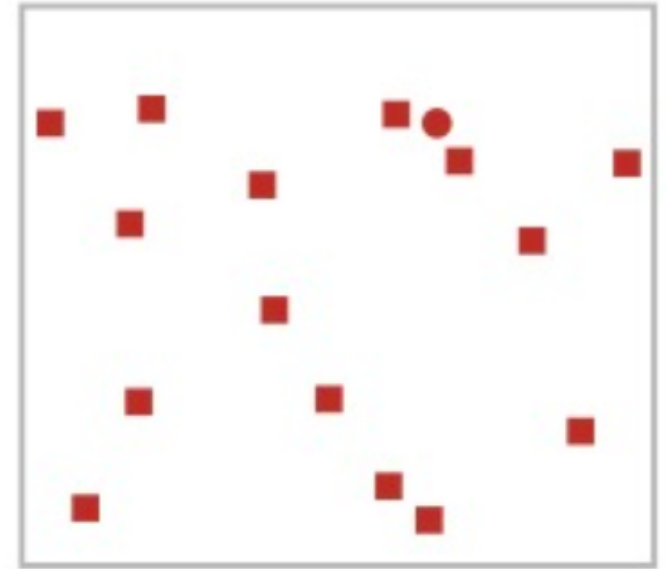
# Popout, or preattentive processing



(a)

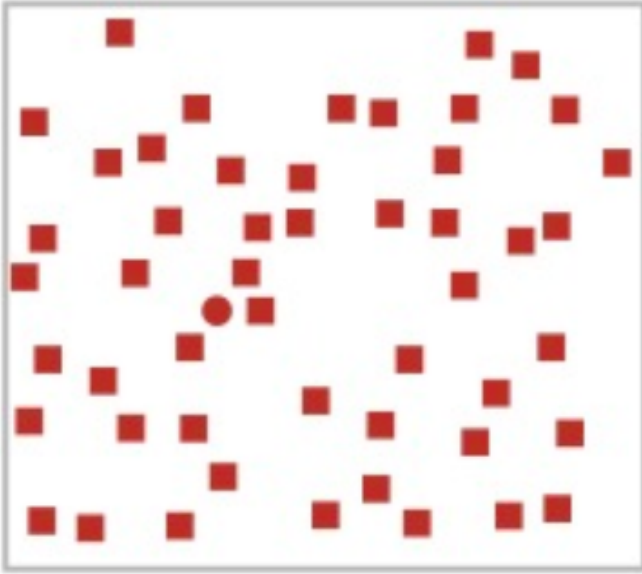


(b)

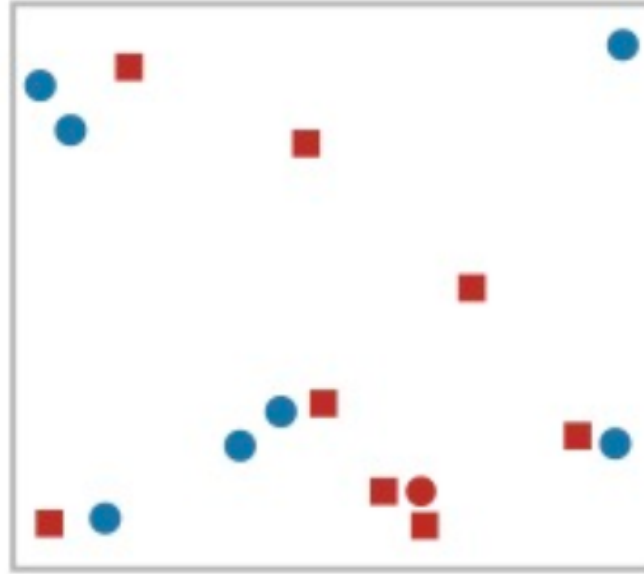


(c)

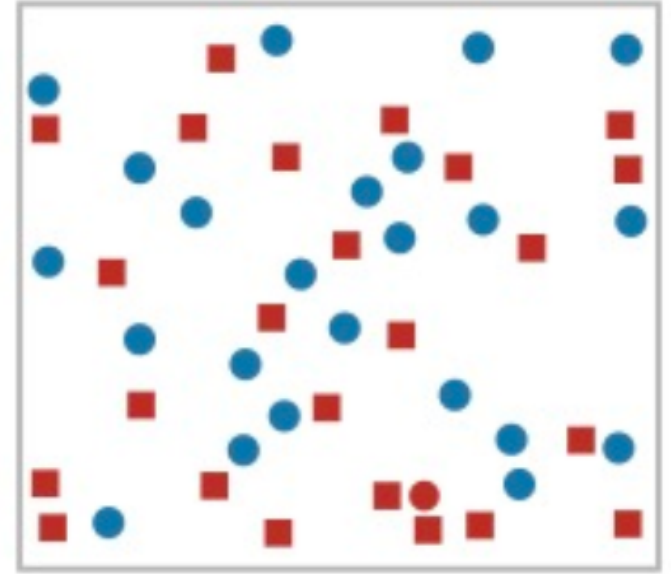
# Popout, or preattentive processing



(d)

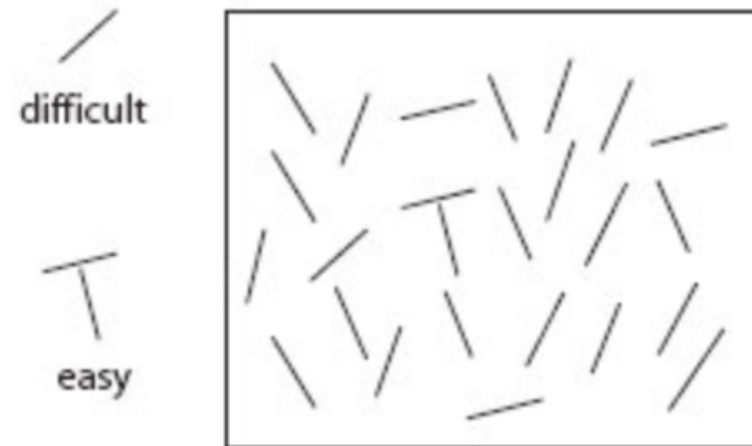
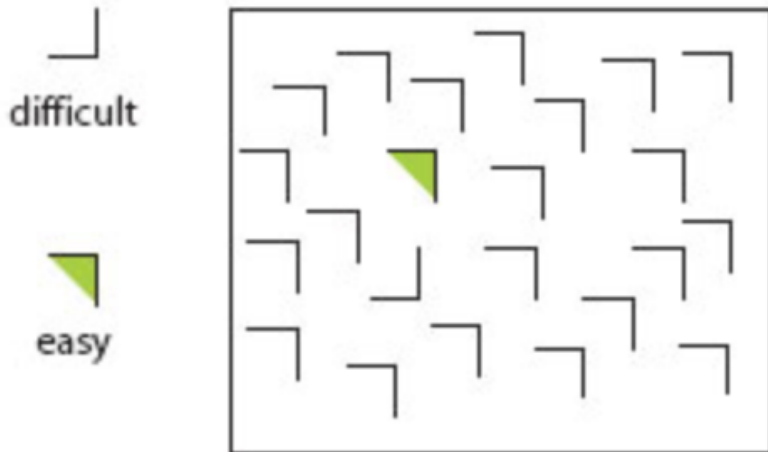
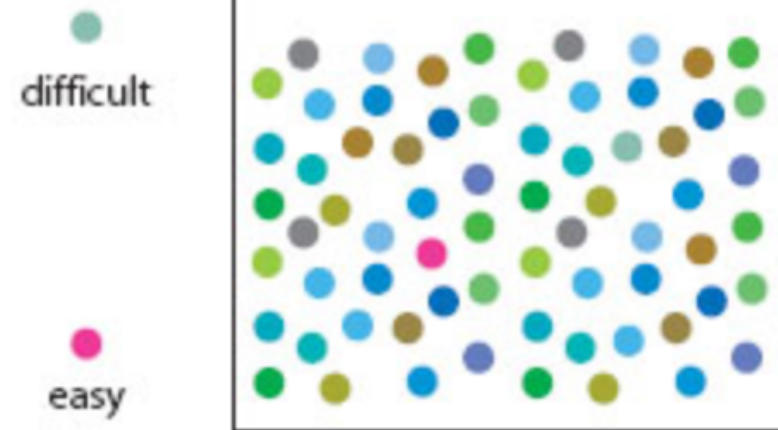
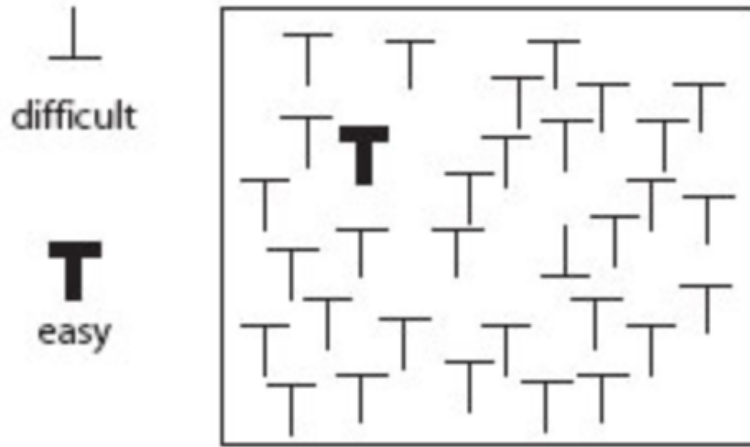


(e)



(f)

# Popout, or preattentive processing

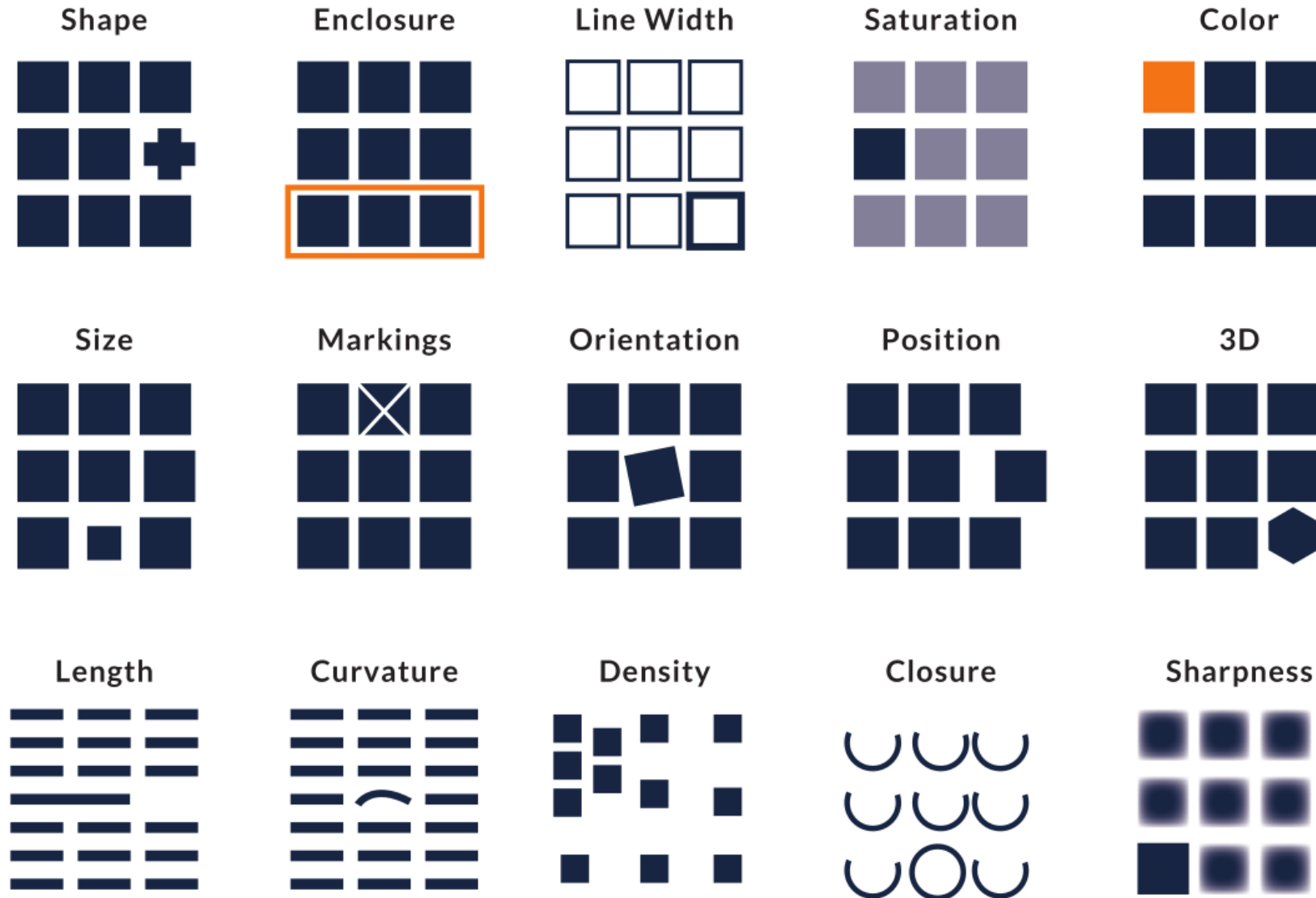




# Popout, or preattentive processing

- The *popout* of the color aesthetic is stronger than the one of the shape aesthetic
- Using many aesthetics at the same time, for different encodings, might weaken the preattentive effect.
- Be very careful at representing several variables all at once using different channels: visual conjunctions are often difficult to see (cfr. separability property)
- In general, if you want something to stand out, make it different from everything else on prominent (possibly combined) visual channels (e.g. color and size)

# Popout, or preattentive processing



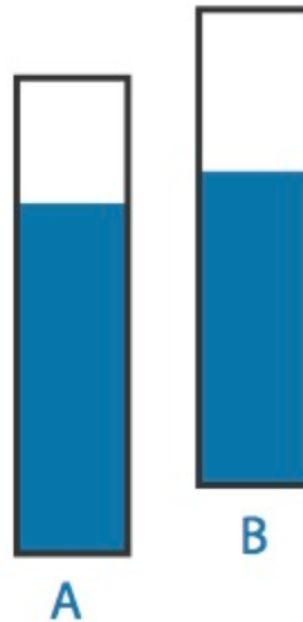
# Relative vs. absolute judgements

Weber's Law: our perceptual system is based on relative judgements, not absolute ones.  
In other terms, our perception is contextual



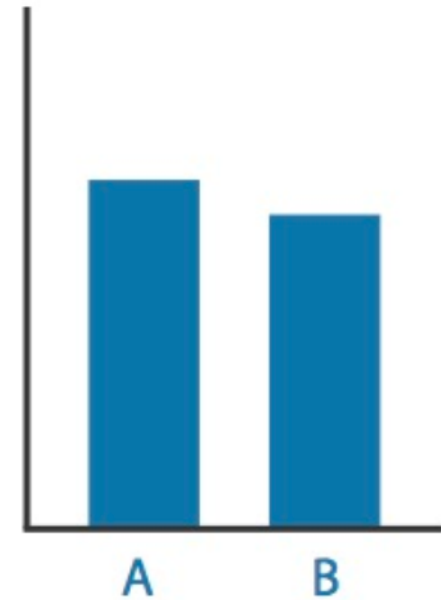
Unframed  
Unaligned

(a)



Framed  
Unaligned

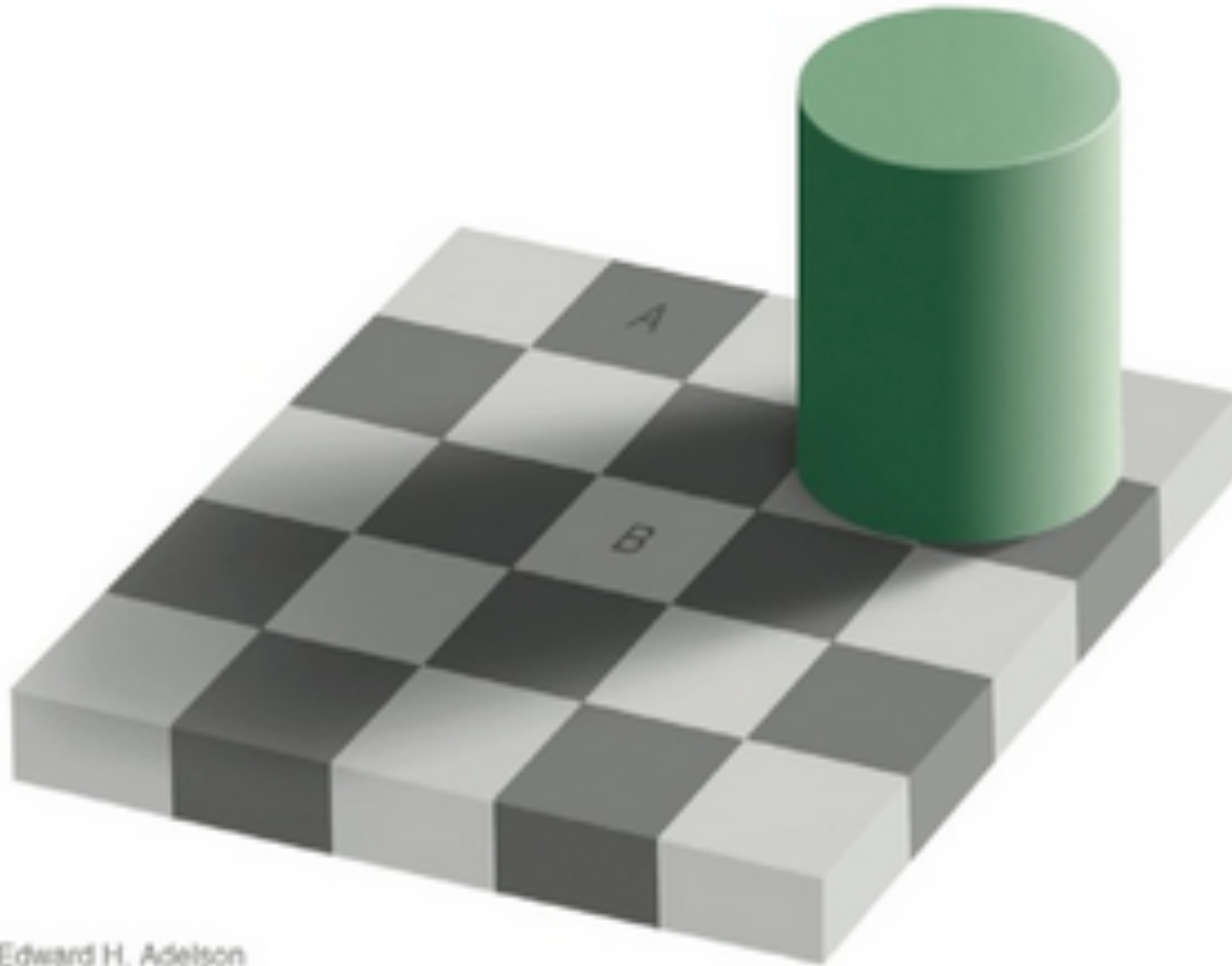
(b)



Unframed  
Aligned

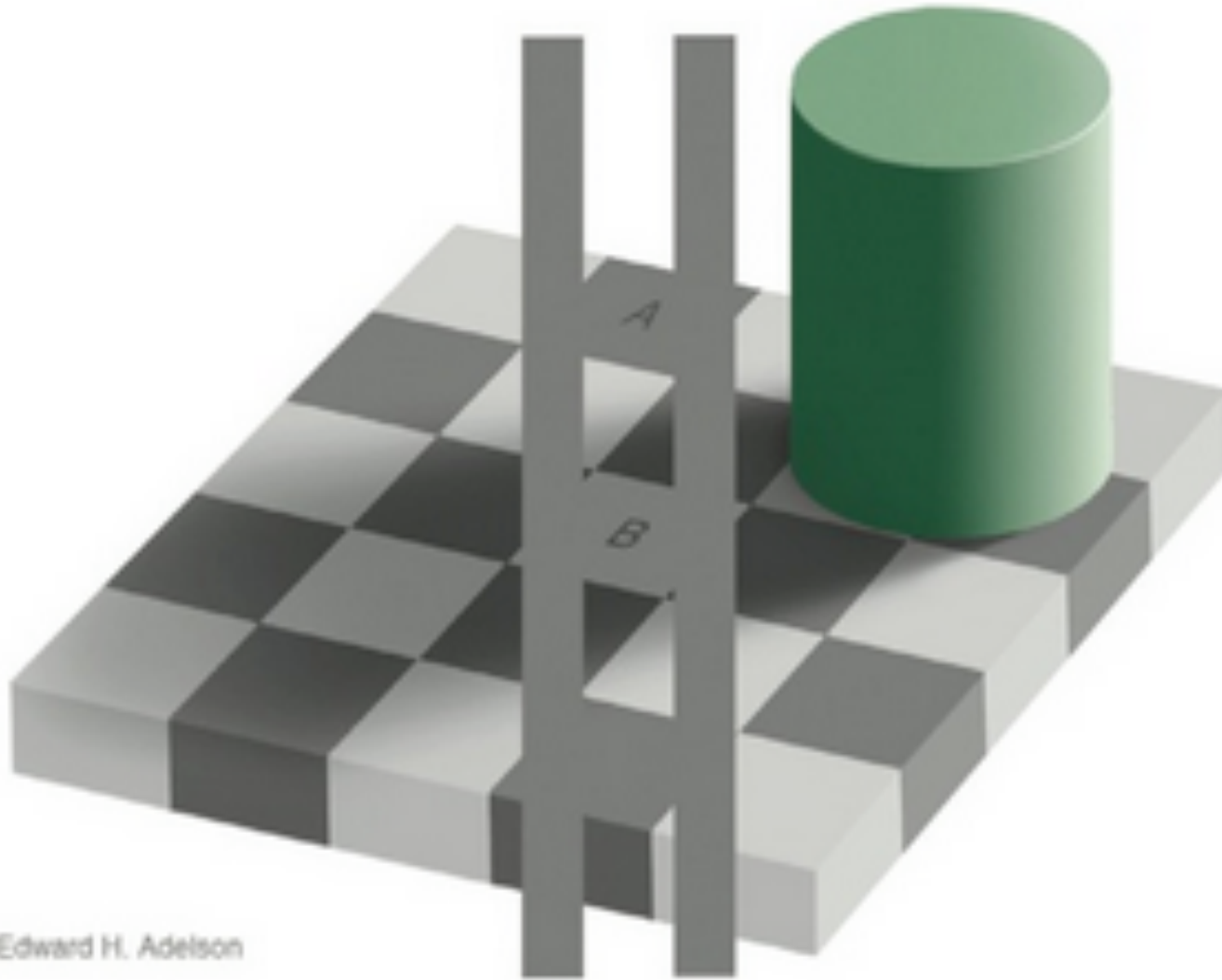
(c)

# Relative vs. absolute judgements



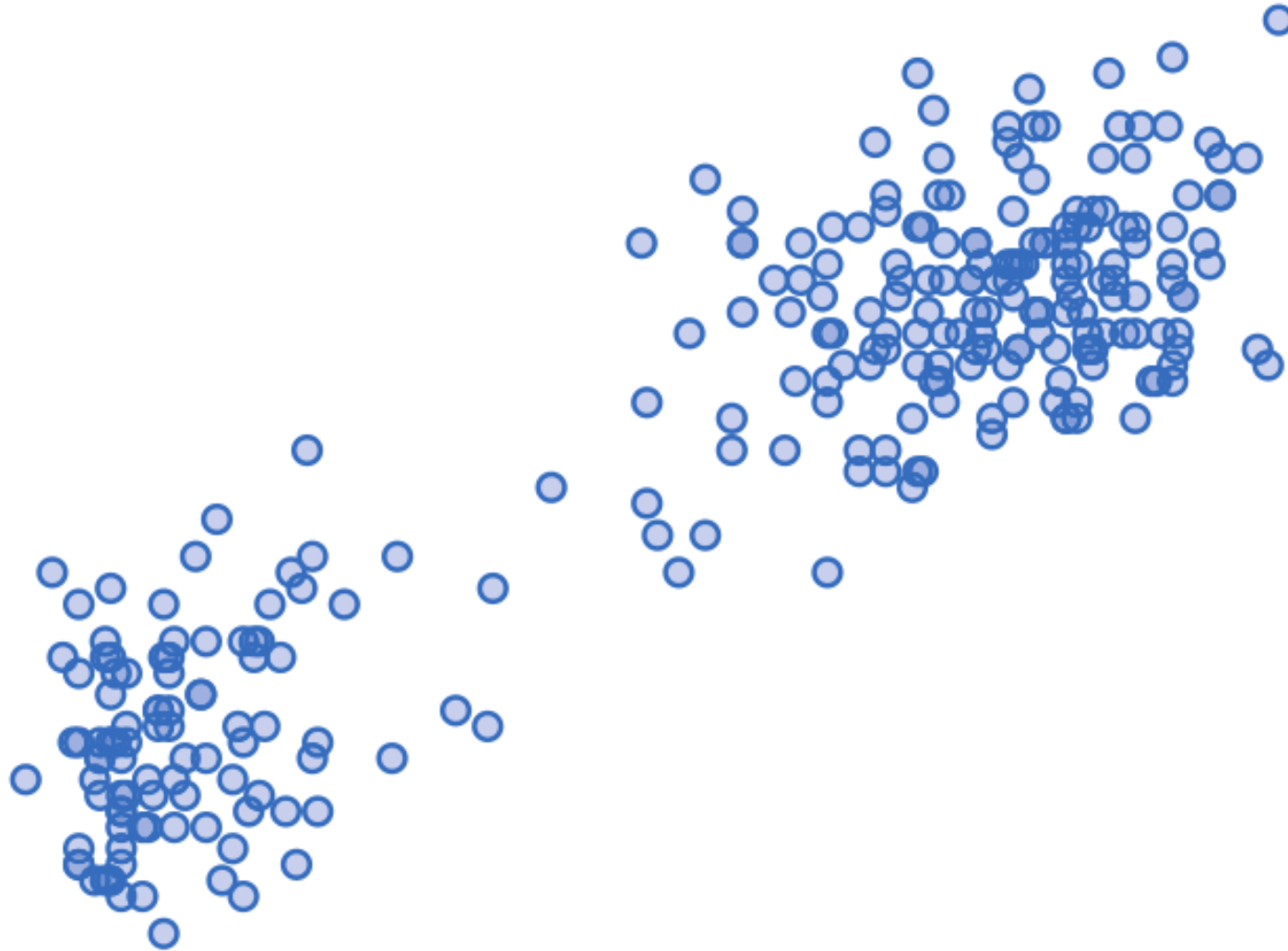
Edward H. Adelson

# Relative vs. absolute judgements

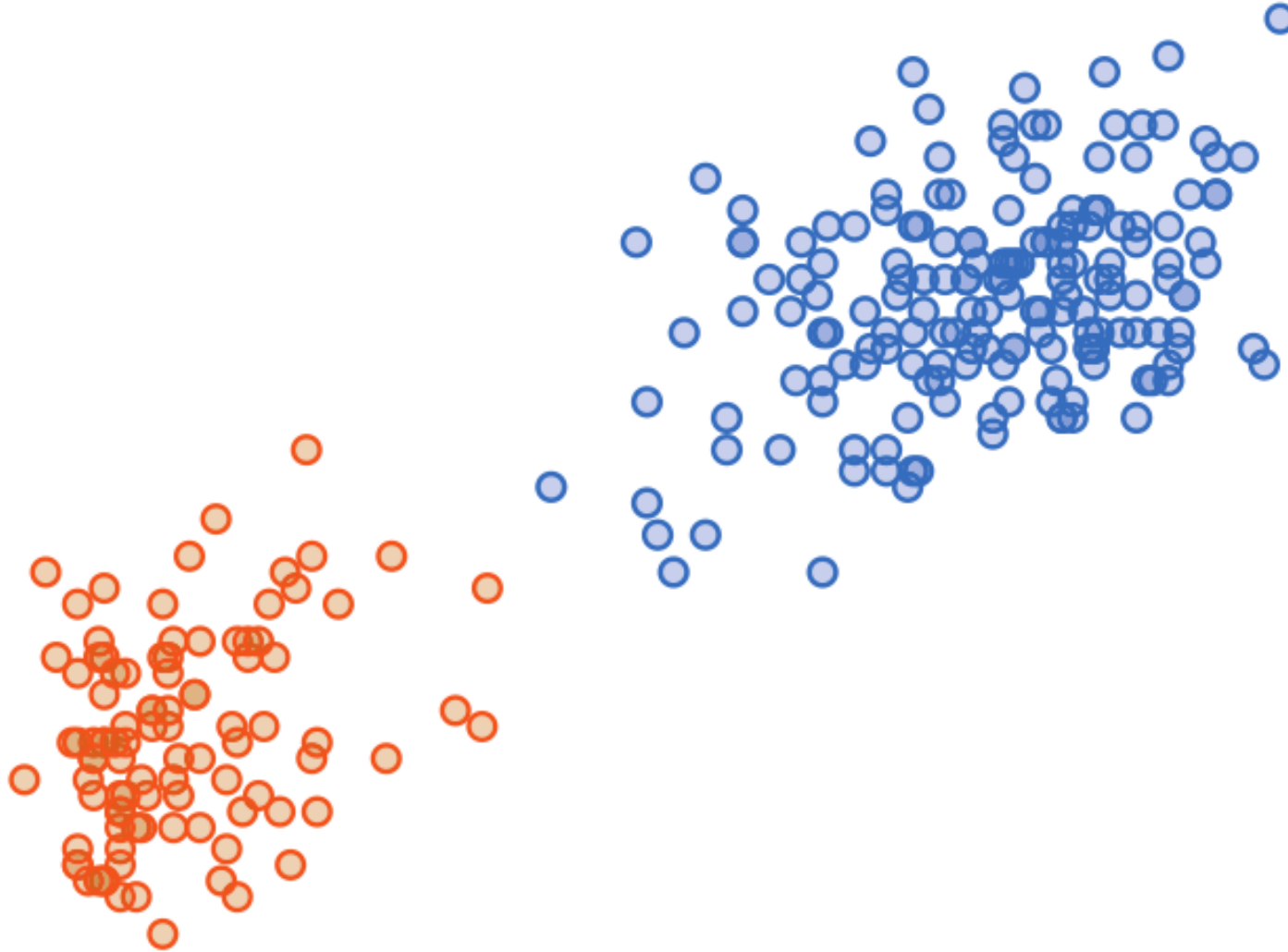


Edward H. Adelson

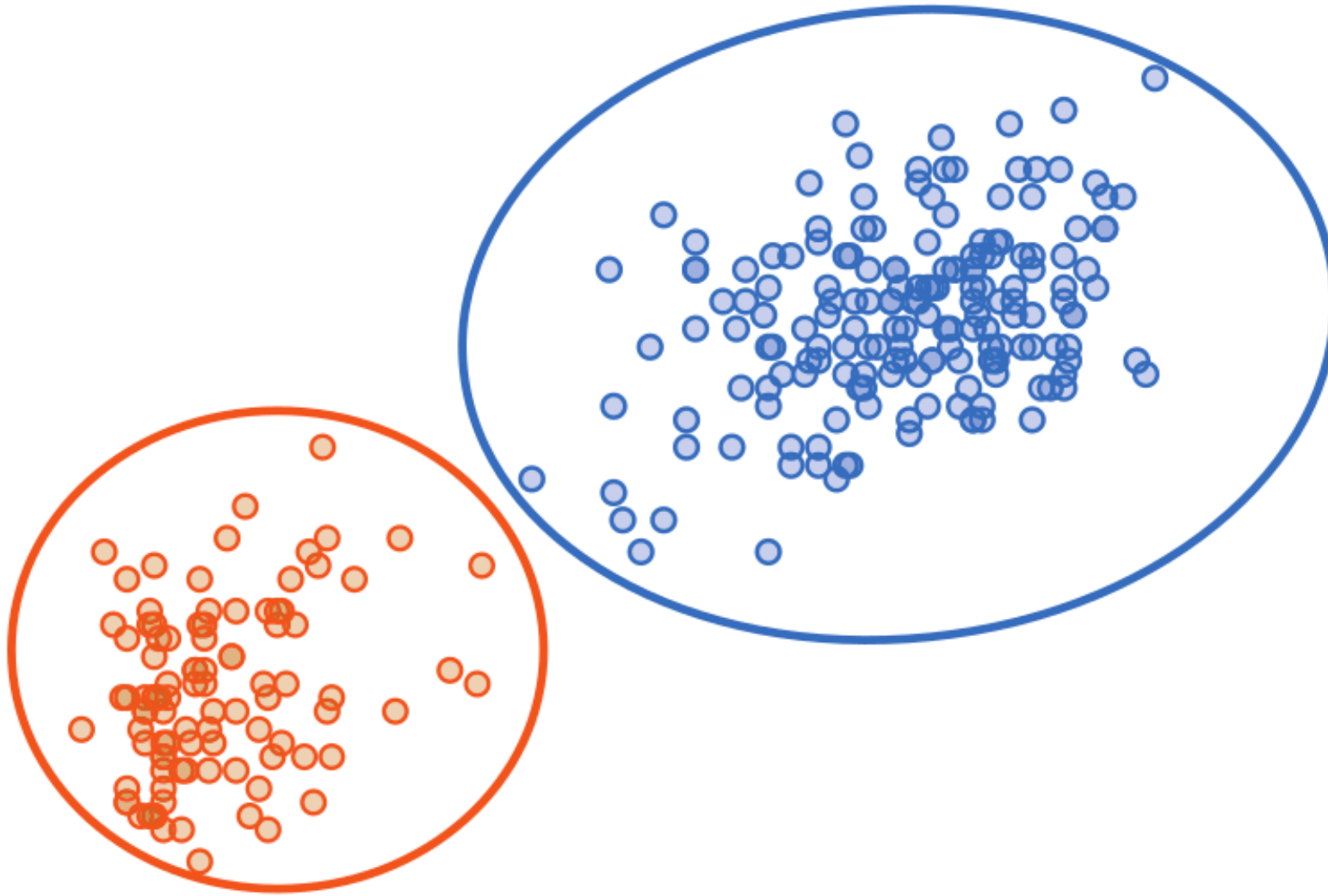
# Gestalt: proximity



# Gestalt: similarity

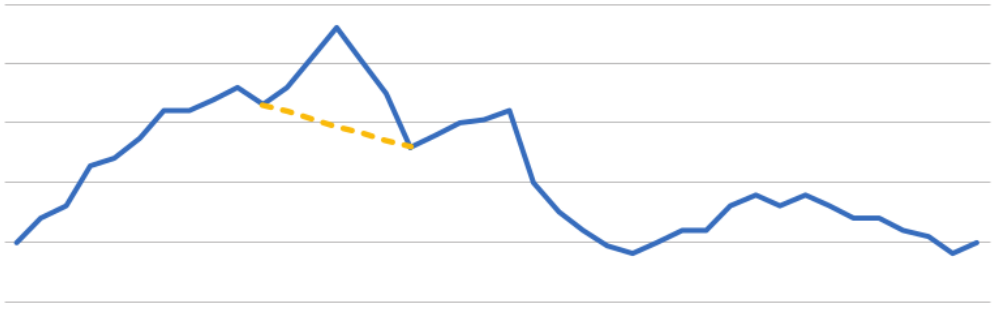


# Gestalt: enclosure

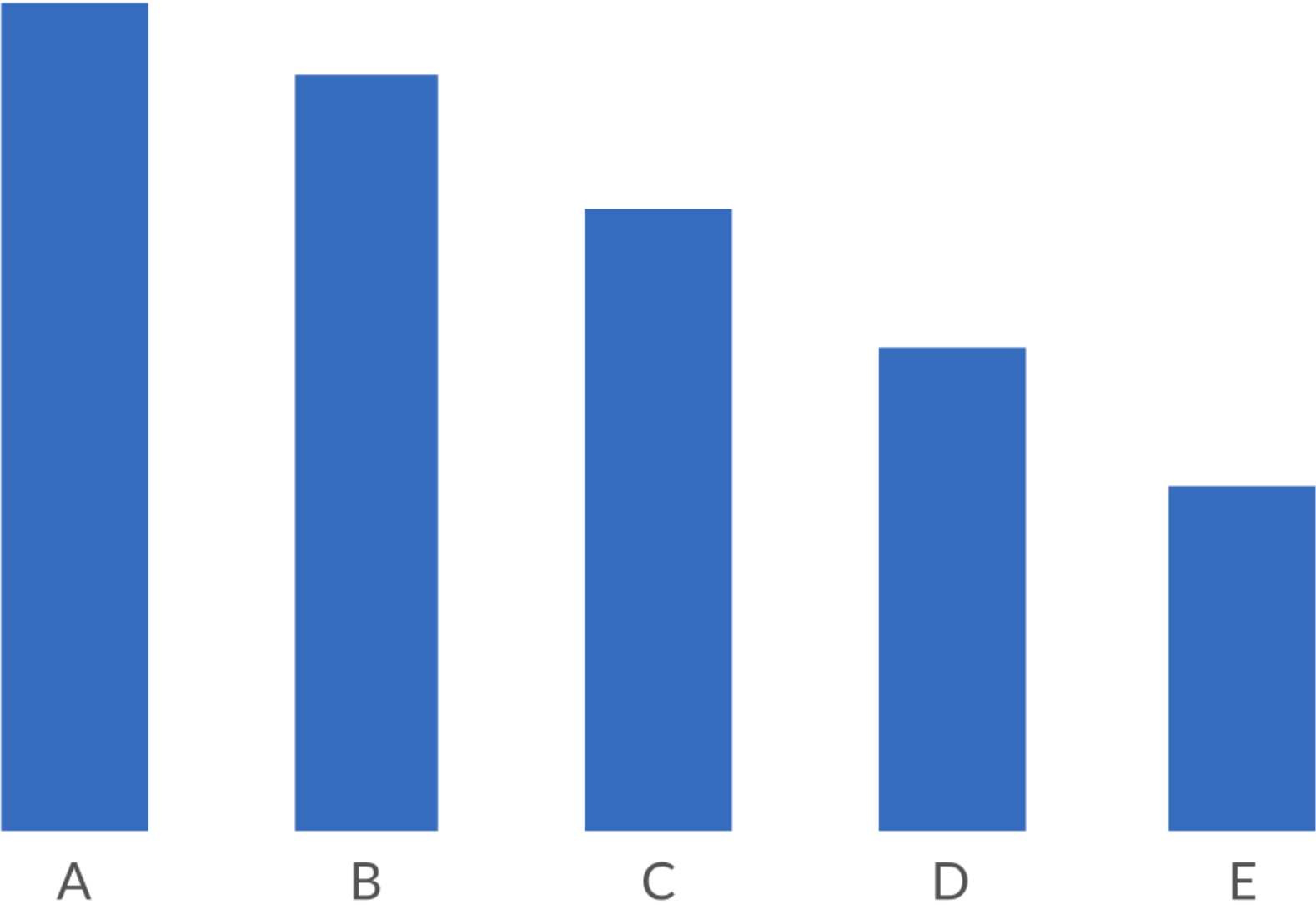




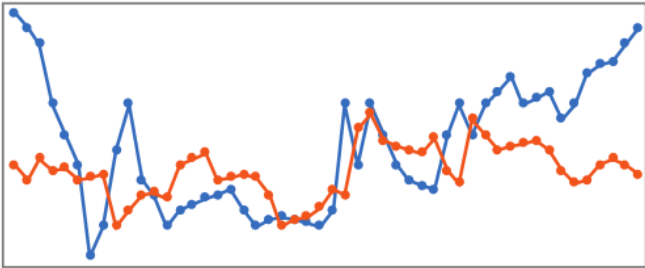
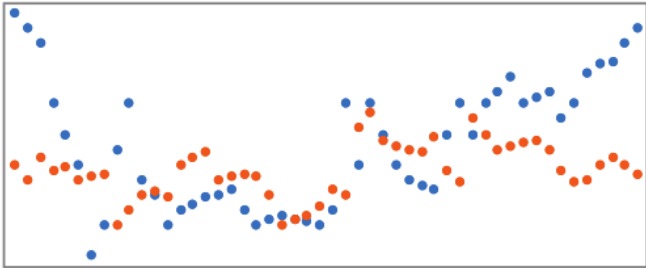
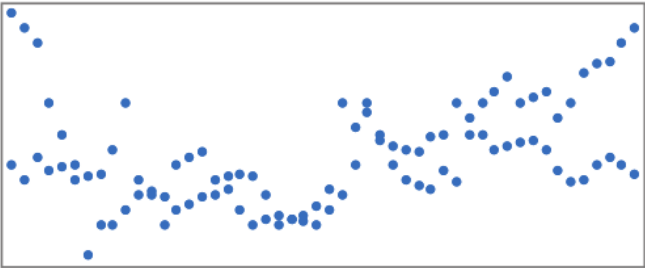
# Gestalt: closure



# Gestalt: continuity



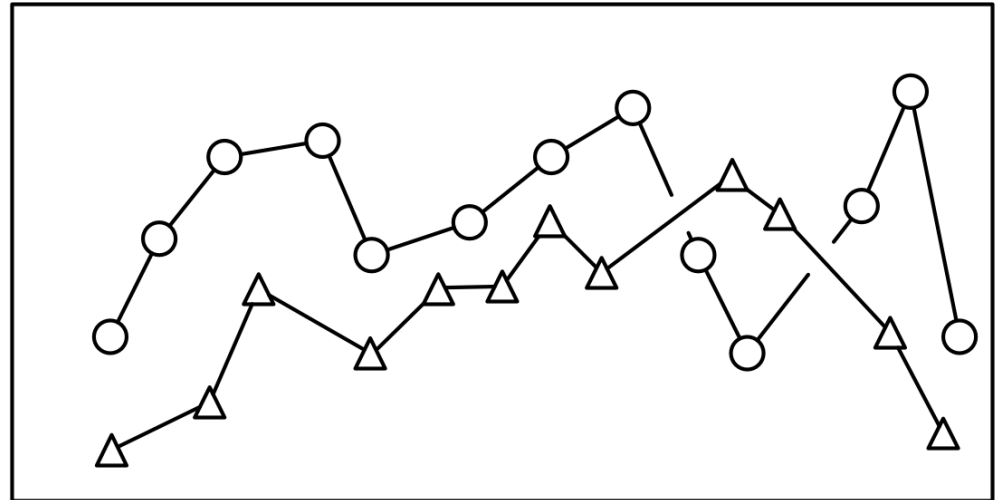
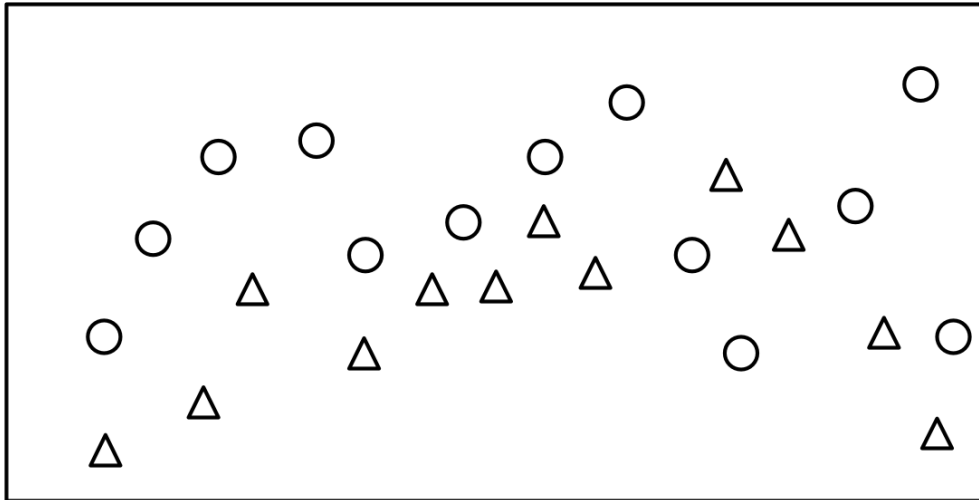
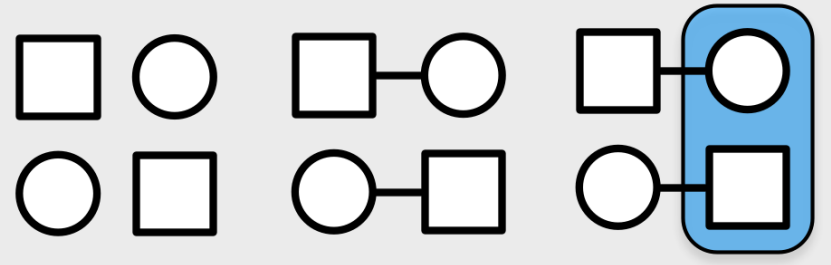
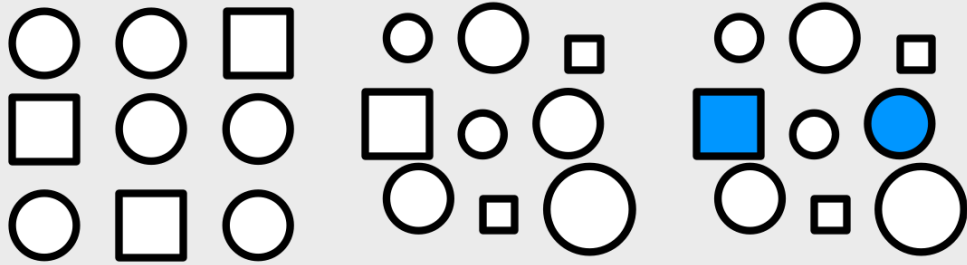
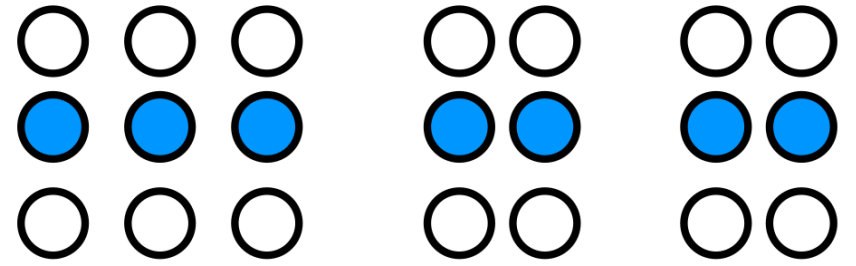
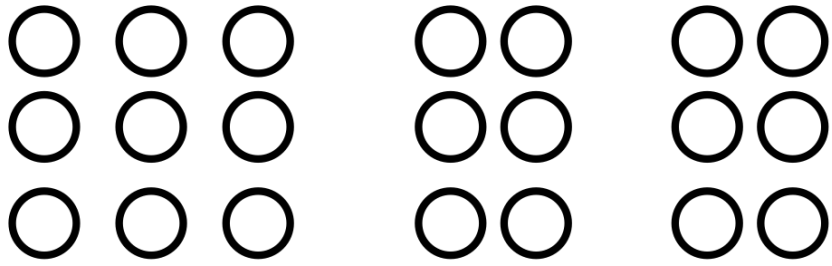
# Gestalt: connection



# Gestalt rules

We look for patterns in what we see:

- Proximity: Things that are spatially near to one another seem to be related.
- Similarity: Things that look alike seem to be related.
- Connection: Things that are visually tied to one another seem to be related.
- Continuity: Partially hidden objects are completed into familiar shapes.
- Closure: Incomplete shapes are perceived as complete.
- Figure and Ground: Visual elements are taken to be either in the foreground or the background.
- Common Fate: Elements sharing a direction of movement are perceived as a unit.



# Reading list

- Munzner: *Visualization Analysis and Design*, chapters 2 and 5 (in the library)
- Ware: *Visual thinking for design*, chapter 2 (in the library)
- Hadley: *R for Data Science*, chapters 15 and 16