# Engineers who make Selfish Machines
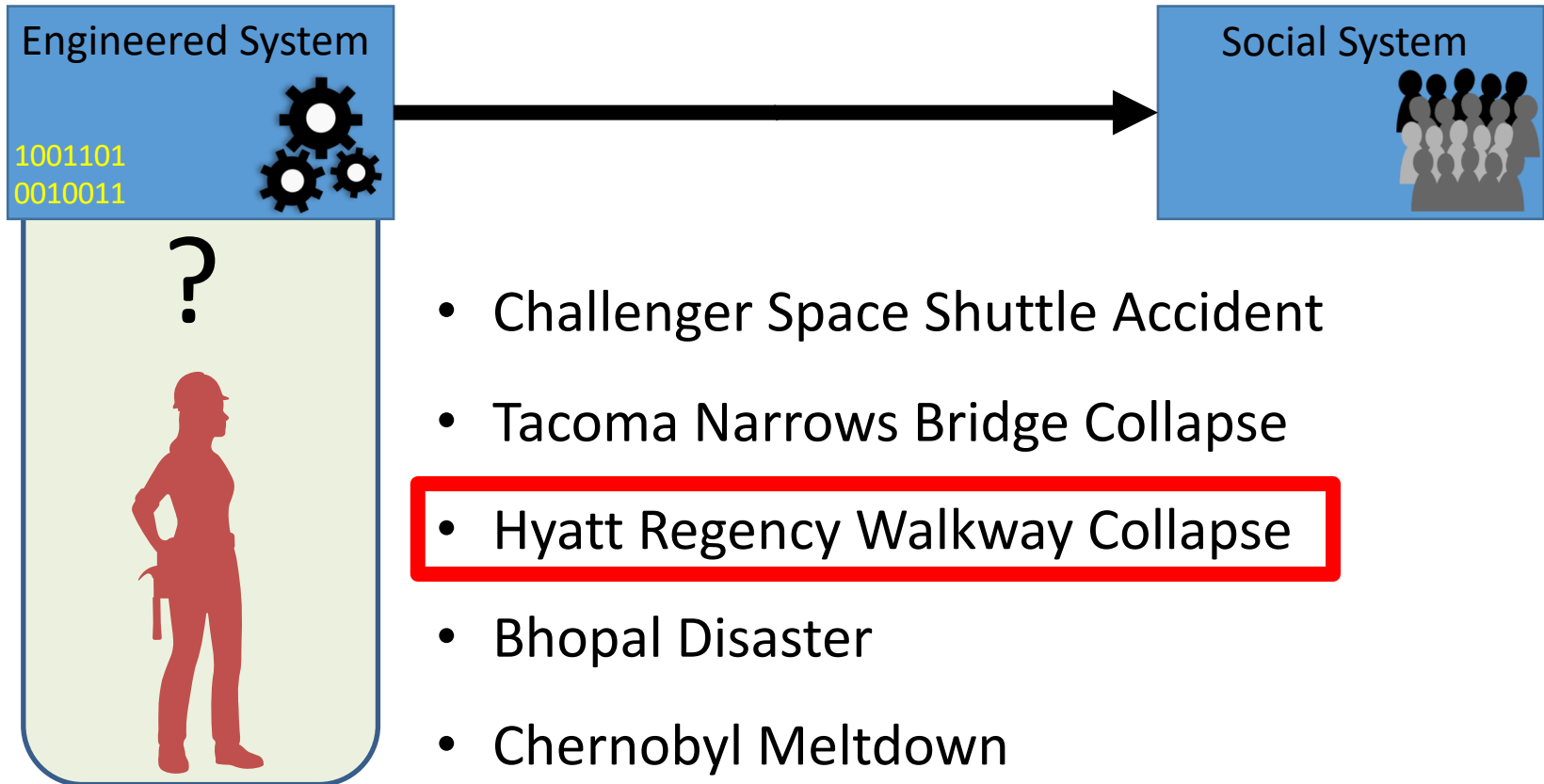
# The Ethics of Socio-Technical Systems

## Philip N. Brown
Dept of Computer Science
University of Colorado, Colorado Springs

University of Colorado
Colorado Springs

University of Colorado
Boulder | Colorado Springs | Denver | Anschutz Medical Campus

Engineered System

1001101
0010011

?

Social System

- Challenger Space Shuttle Accident

- Tacoma Narrows Bridge Collapse

- Hyatt Regency Walkway Collapse

- Bhopal Disaster

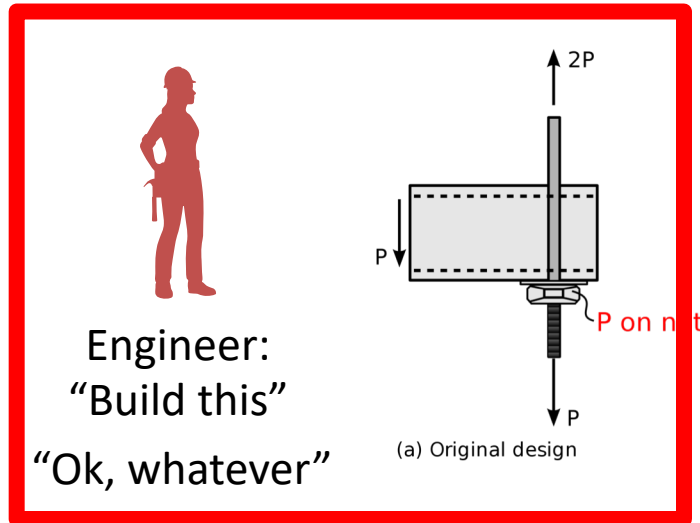- Chernobyl Meltdown

# Hyatt Regency Walkway Collapse

- Kansas City, 1981
- Hotel walkways collapse onto crowded dance floor
- 114 killed, 216 injured
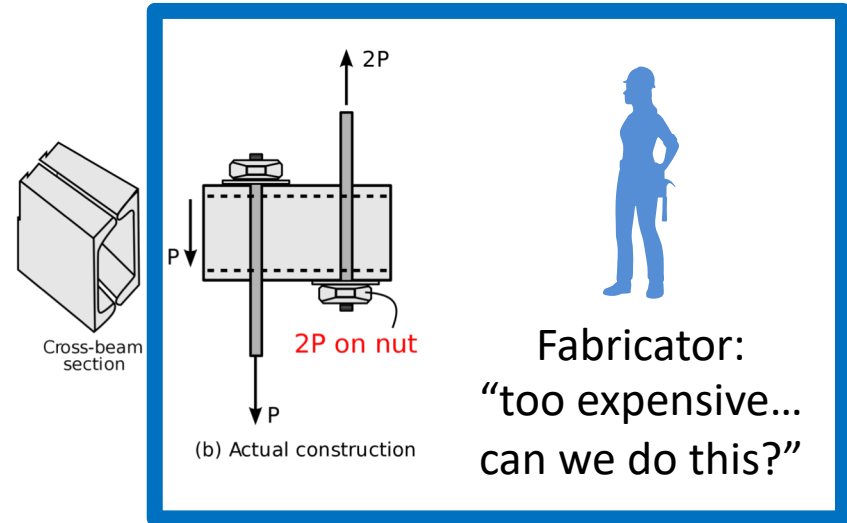- Cause: Systemic neglect of proper design review





Image credit: Public Domain, Dr. Lee Lowery, Jr., P.E.

# Hyatt Regency Walkway Collapse

- Kansas City, 1981
- Hotel walkways collapse onto crowded dance floor
- 114 killed, 216 injured
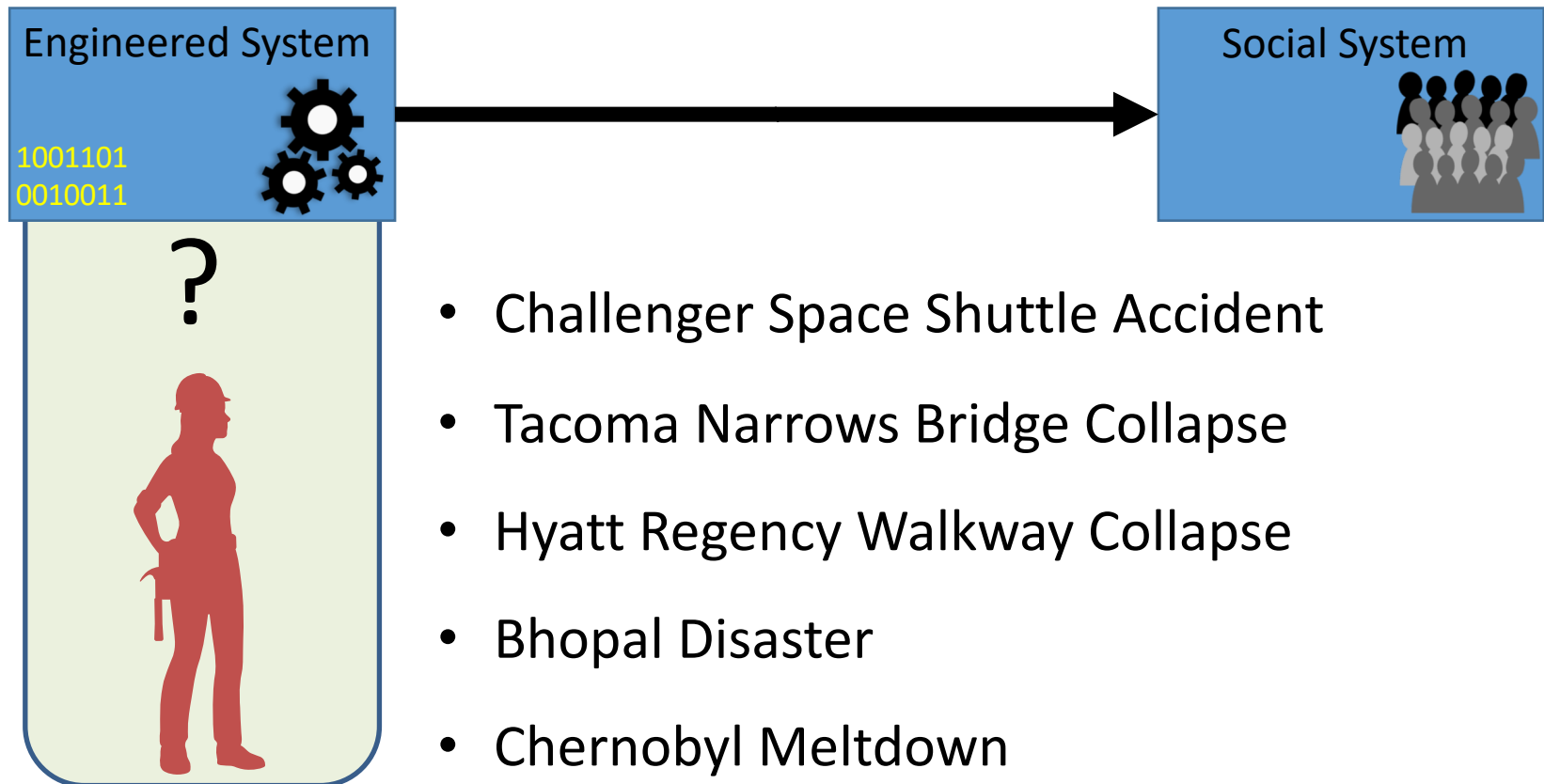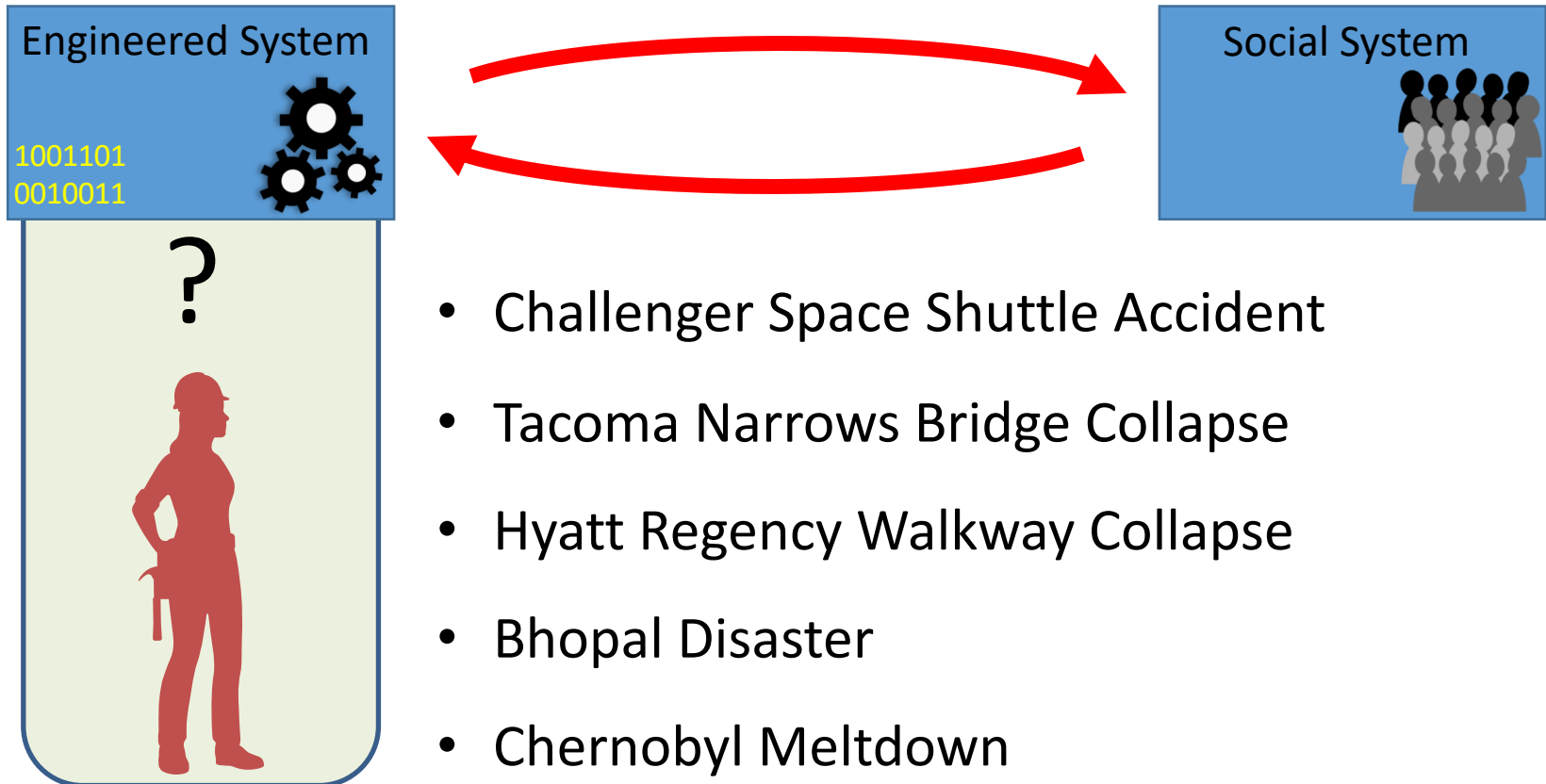- Cause: Systemic neglect of proper design review



Engineer:
"Build this"

"Ok, whatever"

(a) Original design

P on nut

Cross-beam section

2P on nut

(b) Actual construction

Fabricator:
"too expensive…
can we do this?"

Only supports 60% of required load

Only supports **half** load of red design

Image credit: Public Domain, Wikimedia Commons

Engineered System
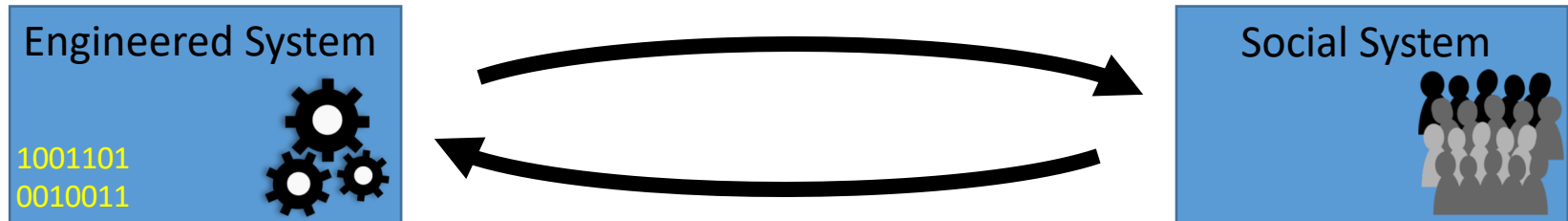
1001101
0010011

?

Social System

- Challenger Space Shuttle Accident

- Tacoma Narrows Bridge Collapse

- Hyatt Regency Walkway Collapse

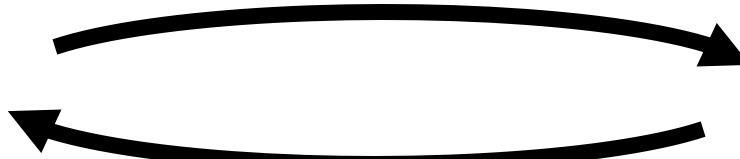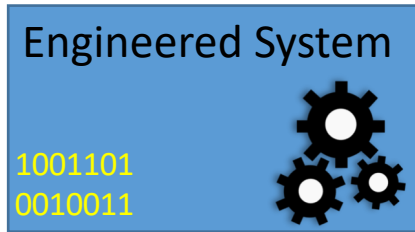- Bhopal Disaster

- Chernobyl Meltdown

**Ethics Message: Engineer things that don't break**

**Implication: Adequate design is sufficient.**

Engineered System

1001101
0010011

?
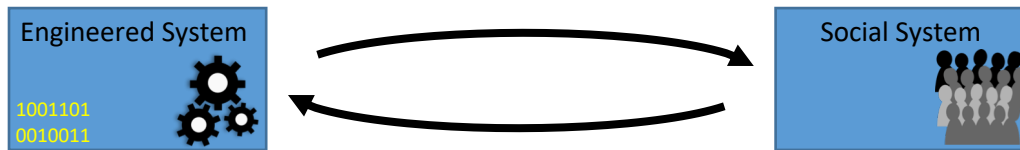
Social System

- Challenger Space Shuttle Accident

- Tacoma Narrows Bridge Collapse

- Hyatt Regency Walkway Collapse

- Bhopal Disaster

- Chernobyl Meltdown

**Ethics Message: Engineer things that don't break**

**Implication: Adequate design is sufficient.**

Engineered System

1001101
0010011

Social System

- Why does the feedback loop matter?

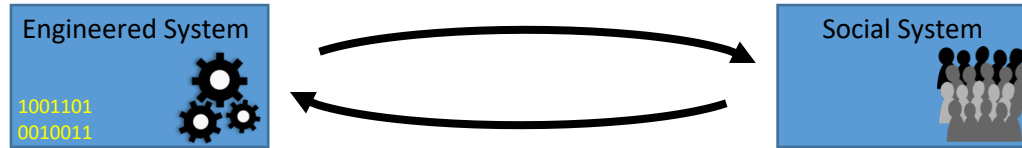- What are ethical implications?

- How to teach this?

Engineered System

1001101
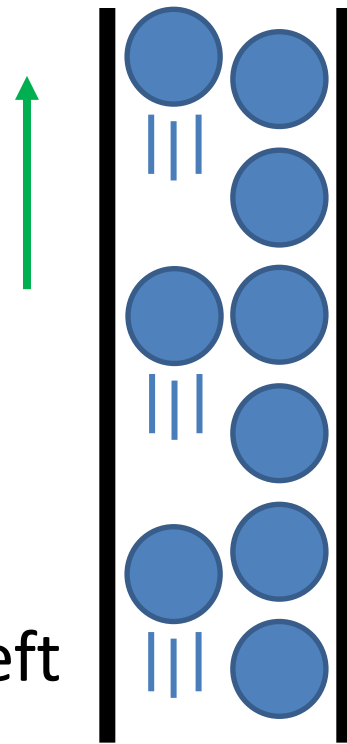0010011

Social System

← Bat Cave

Concourse A
Terminal
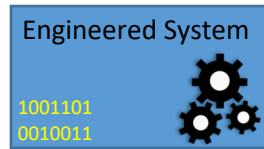Luggage Claim
Hotel

Stand on right

Image credit: Public Domain, Pixabay

University of Colorado
Colorado Springs

University of Colorado
Boulder | Colorado Springs | Denver | Anschutz Medical Campus

Walk on left | Stand on right

Image credit: Public Domain, Pixabay

University of Colorado
Colorado Springs

University of Colorado
Boulder | Colorado Springs | Denver | Anschutz Medical Campus

Engineered System

1001101
0010011

Social System

Walk on left

Stand on right

Engineered System

1001101
0010011

Social System

Human Intuition:
Don't be that guy!

Walk on left

Stand on right

University of Colorado
Colorado Springs

University of Colorado
Boulder | Colorado Springs | Denver | Anschutz Medical Campus

Engineered System

1001101
0010011

Social System

Human Intuition:
Don't be that guy!

Walk on left

Stand on right

Question: Robot's escalator policy?

University of Colorado
Colorado Springs

University of Colorado
Boulder | Colorado Springs | Denver | Anschutz Medical Campus

Low Impact

Selfish

Stand Left

Question: Robot's escalator policy?

Low Impact

Selfish

Stand Left

Human Intuition: not this
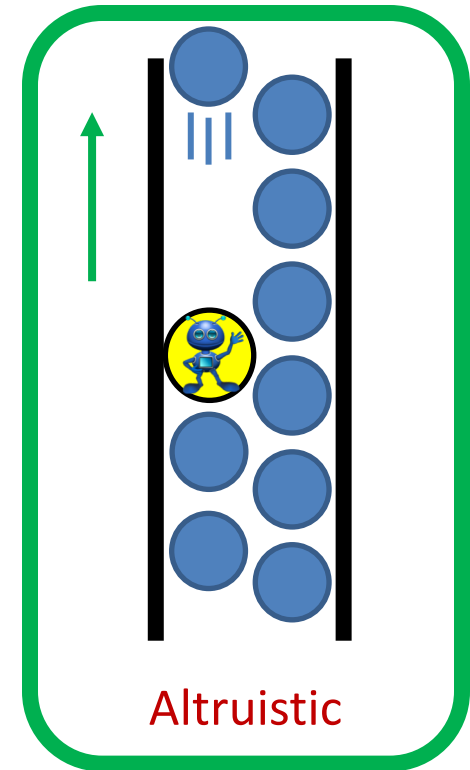
Question: Robot's escalator policy?
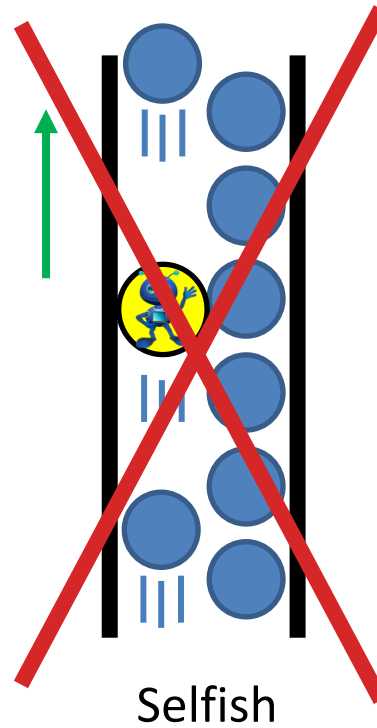
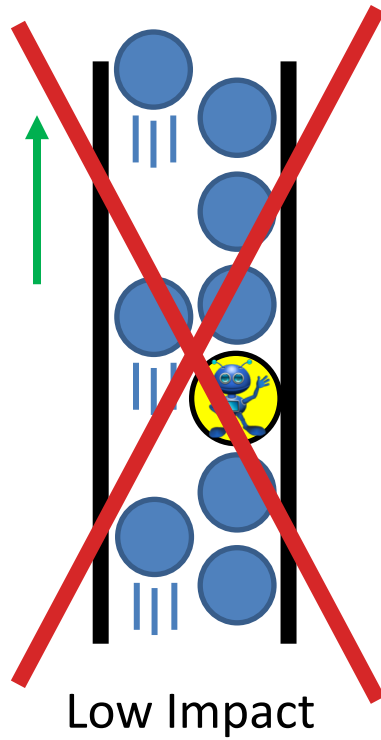Low Impact

Selfish

Stand Left

Systems Perspective:
definitely this
(maximizes throughput)

For more info: Prof. Lesley Strawderman at Mississippi State University

Low Impact

Selfish
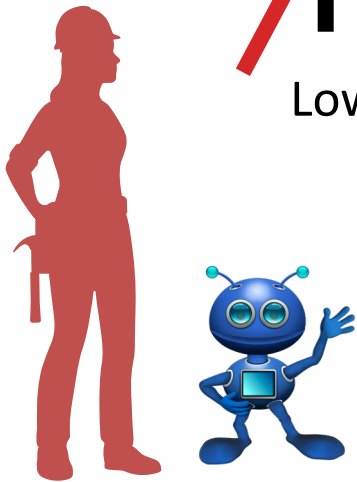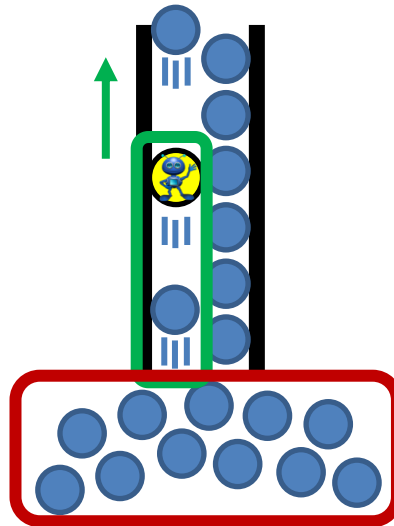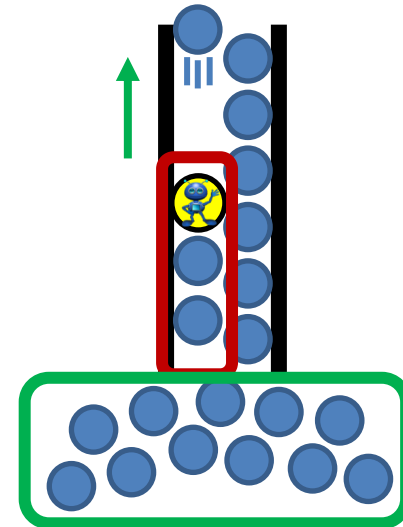
Altruistic

Systems Perspective:
definitely this
(maximizes throughput)

For more info: Prof. Lesley Strawderman at Mississippi State University

University of Colorado
Colorado Springs

University of Colorado
Boulder | Colorado Springs | Denver | Anschutz Medical Campus

**Selfish**

**Altruistic**



**Benefits:** Robot / a few people behind — Everybody waiting / to use escalator

**Harms:** Everybody waiting / to use escalator — Robot / a few people behind

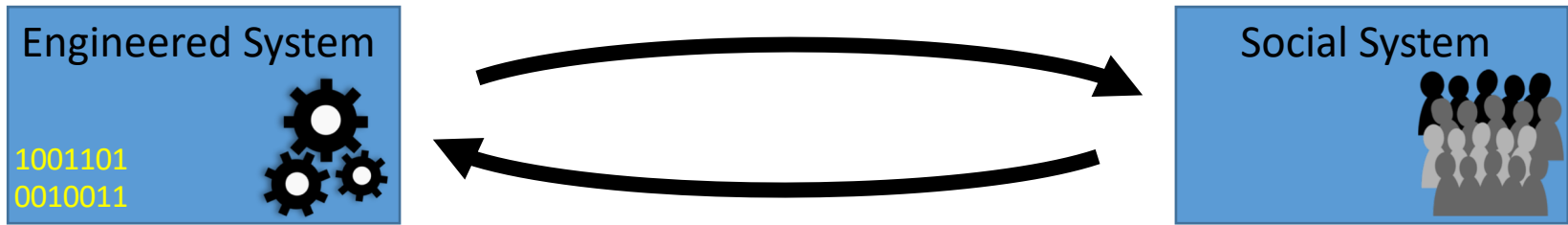**Ethical?** Bad for "system" / Robot acts like a human — Good for "system" / Robot looks like a jerk!

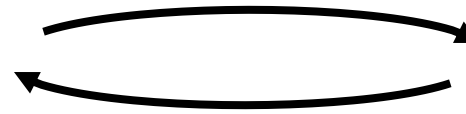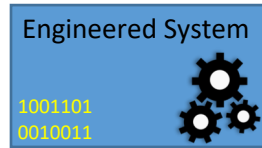For more info: Prof. Lesley Strawderman at Mississippi State University

Engineered System

1001101
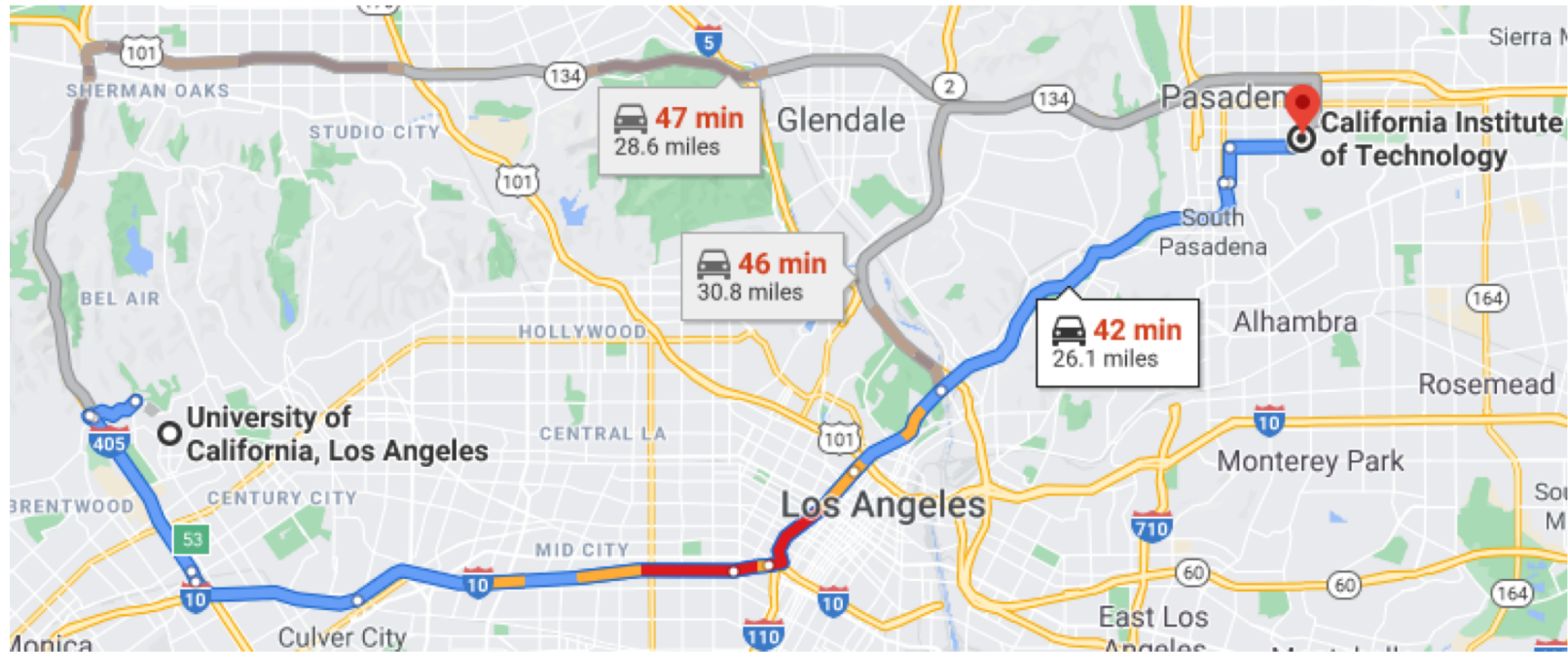0010011

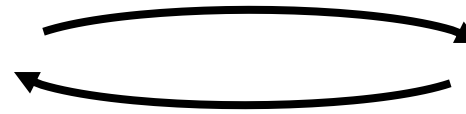Social System
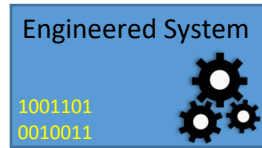
**Traditional Ethics Message:** Engineer things that don't break

Engineered System

1001101
0010011

Social System

**Needed Update:** Engineer your machines to *interact* with people

University of Colorado
Colorado Springs

University of Colorado
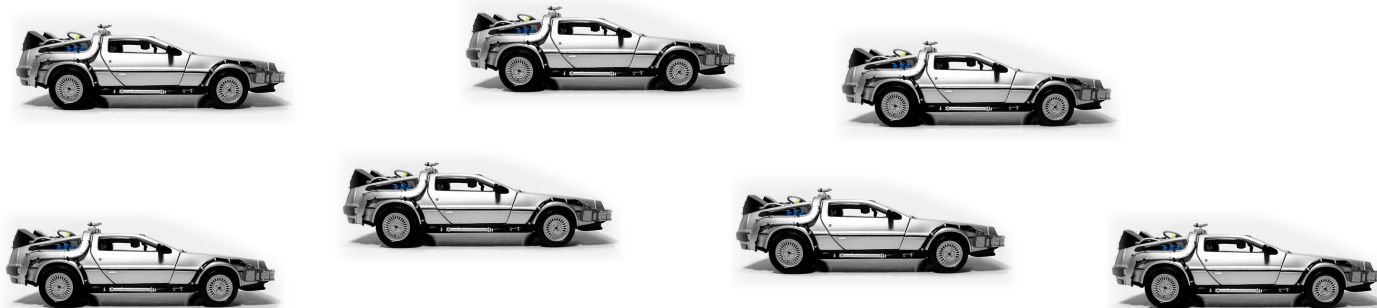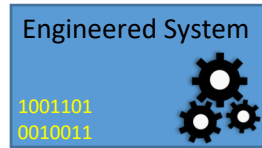Boulder | Colorado Springs | Denver | Anschutz Medical Campus

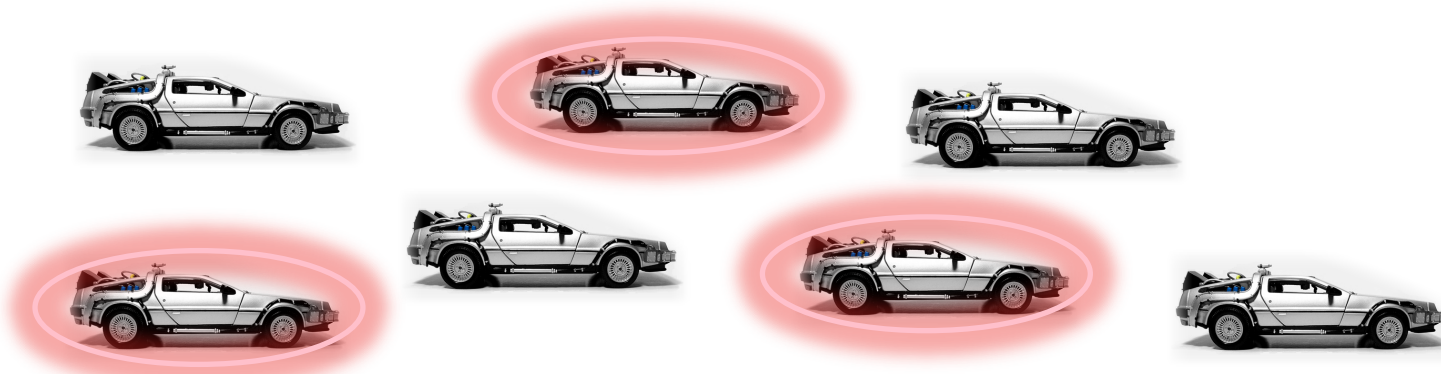# Choosing Routes in Highway Networks

# Choosing Routes in Highway Networks

Agenda: pose simple model
Explore ethics in its context

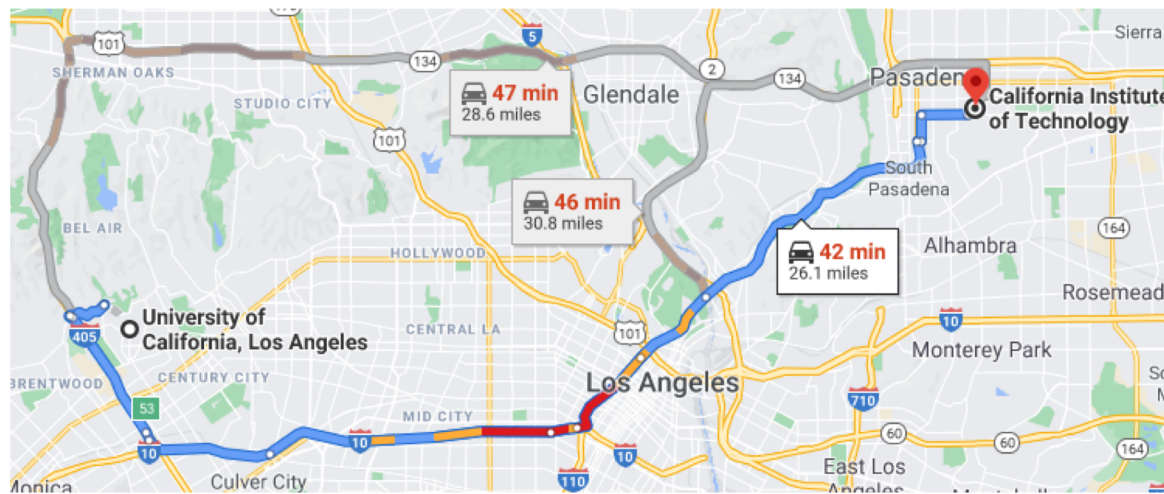Engineered System

1001101
0010011

Social System

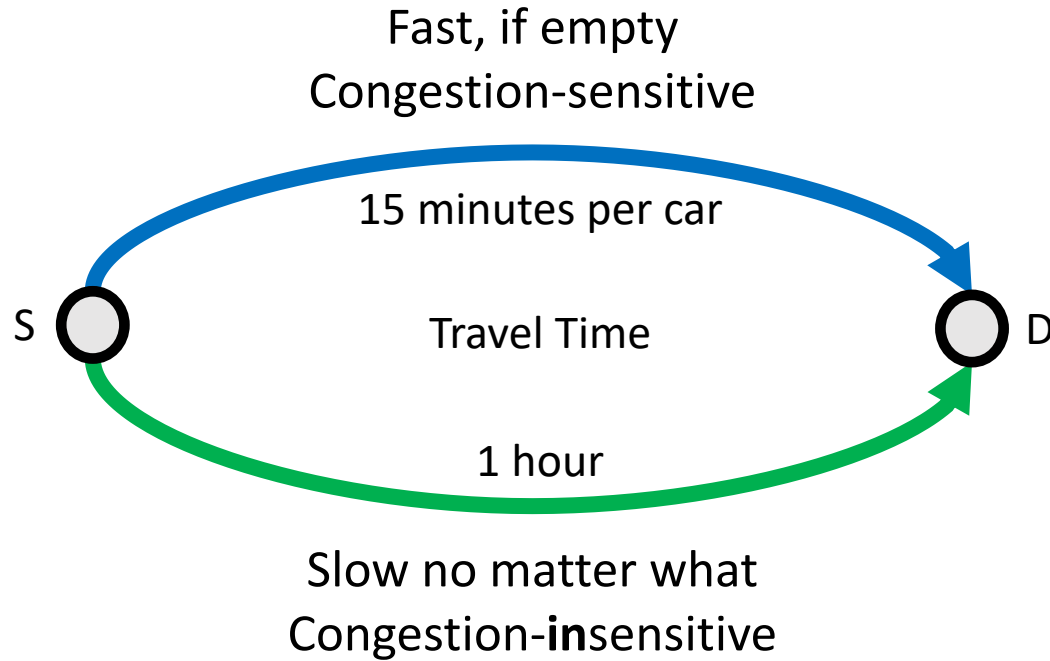# Choosing Routes in Highway Networks

?
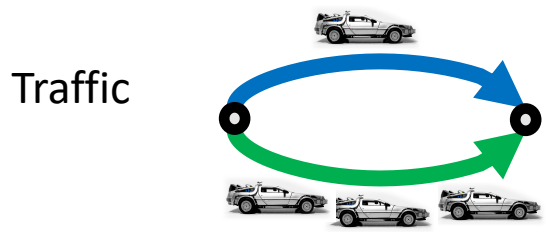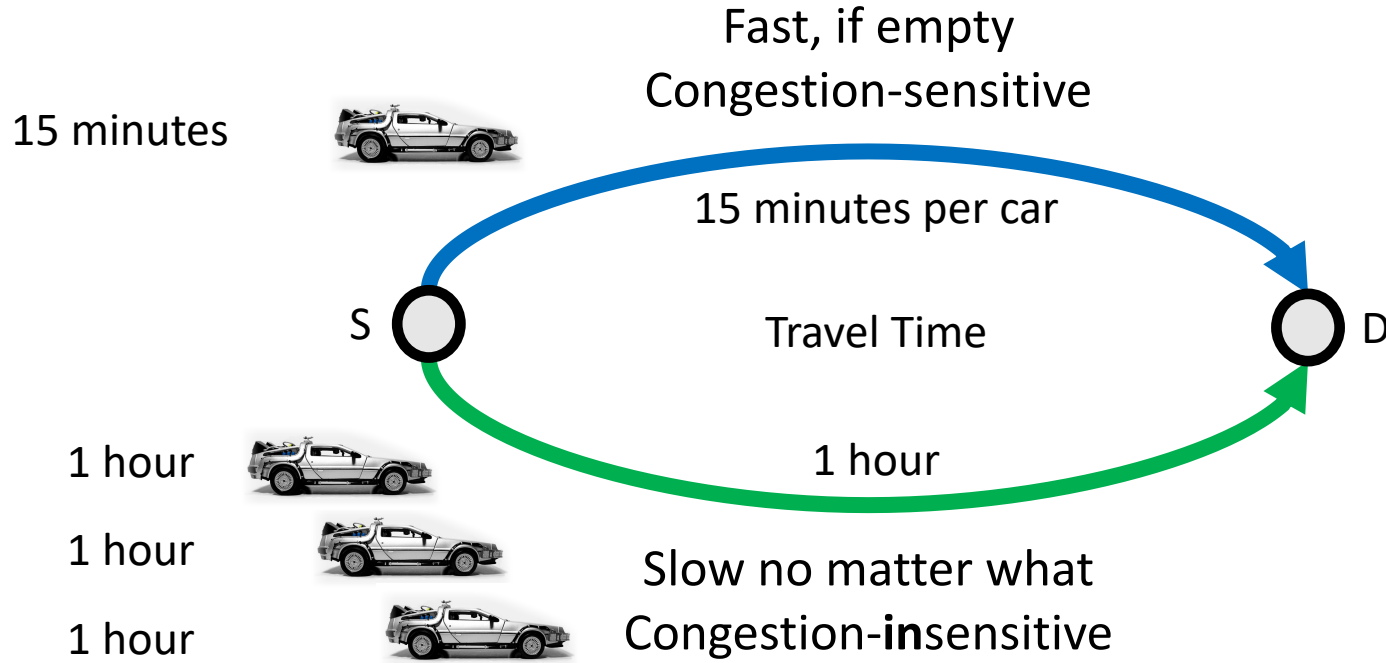
Agenda: pose simple model
Explore ethics in its context

Question: should self-driving cars be altruistic?

Fast, if empty
Congestion-sensitive

15 minutes per car

S          Travel Time          D

1 hour

Slow no matter what
Congestion-**in**sensitive

Fast, if empty
Congestion-sensitive

15 minutes

15 minutes per car

S

Travel Time

D

1 hour

1 hour

1 hour

1 hour

Slow no matter what
Congestion-**in**sensitive
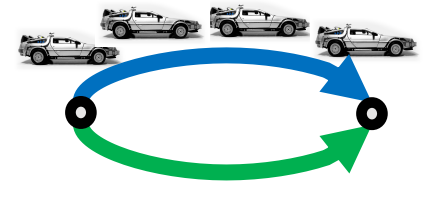
Traffic

Total
Time

3.25 hours

30 minutes

30 minutes

Fast, if empty
Congestion-sensitive

15 minutes per car

S

Travel Time

D

1 hour

1 hour

1 hour

Slow no matter what
Congestion-**in**sensitive

Traffic

Total
Time

3.25 hours

3 hours

45 minutes

45 minutes

45 minutes

Fast, if empty
Congestion-sensitive

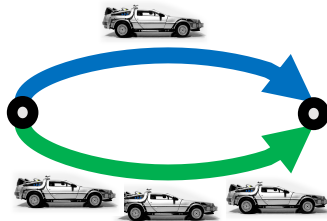15 minutes per car

Travel Time

S          D

1 hour

1 hour

Slow no matter what
Congestion-**in**sensitive

Traffic

Total
Time

3.25 hours          3 hours          3.25 hours

1 hour

1 hour

1 hour

1 hour
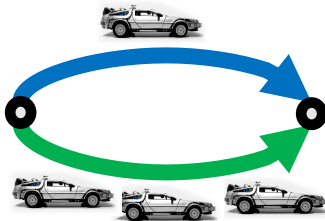
Fast, if empty
Congestion-sensitive

15 minutes per car

Travel Time

S

D

1 hour

Slow no matter what
Congestion-**in**sensitive

Traffic

Total
Time

3.25 hours          3 hours          3.25 hours          4 hours

University of Colorado
Colorado Springs

University of Colorado
Boulder | Colorado Springs | Denver | Anschutz Medical Campus

1 hour

1 hour

1 hour

1 hour

15 minutes per car

S

Travel Time

1 hour

D

Best option!
(Pareto optimal)

Traffic

Total
Time

3.25 hours

3 hours

3.25 hours

4 hours

30 minutes

30 minutes

15 minutes per car

S — D

Travel Time

1 hour
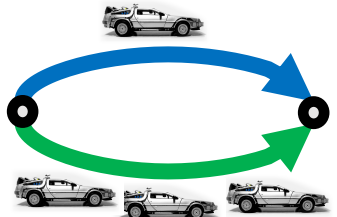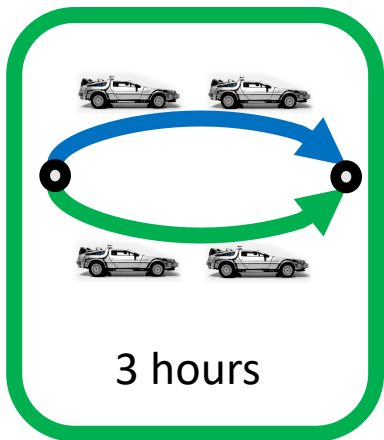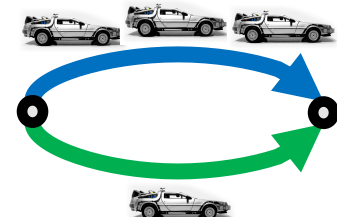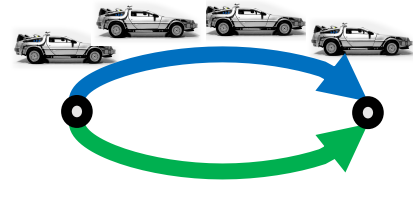
1 hour

1 hour

Best option!
(Pareto optimal)

Traffic

Total
Time

3.25 hours    3 hours    3.25 hours    4 hours

30 minutes

30 minutes

**Altruistic self-driving cars?**

15 minutes per car

S
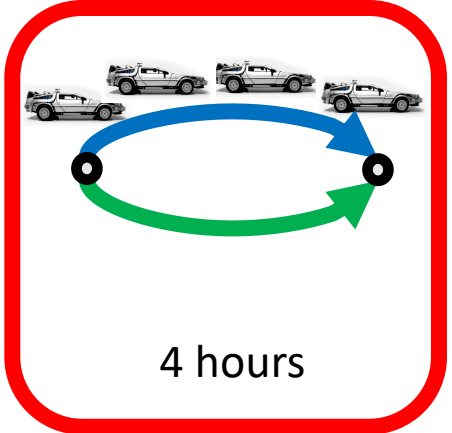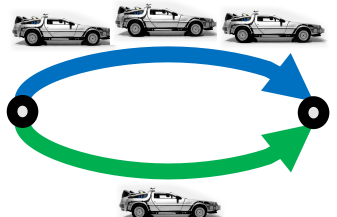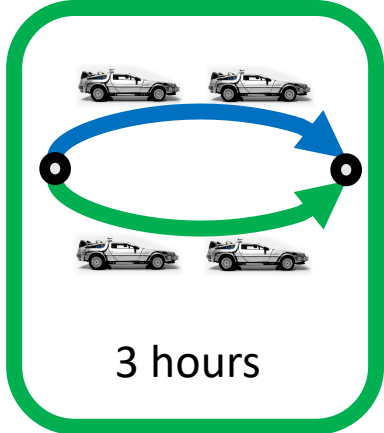
Incentive to switch!
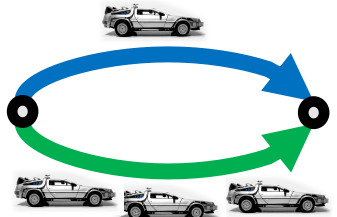
D

1 hour

1 hour

1 hour
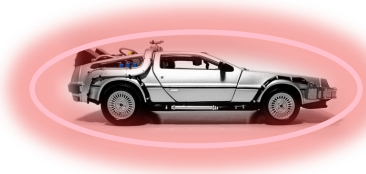
Best option!
(Pareto optimal)

Selfish traffic
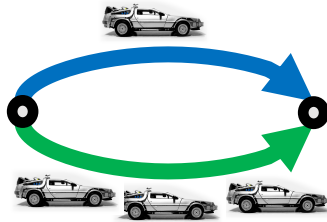is like this!

Traffic

Total
Time

3.25 hours

3 hours

3.25 hours

4 hours

# Altruistic self-driving cars?

Altruism: act like there is 2x actual traffic

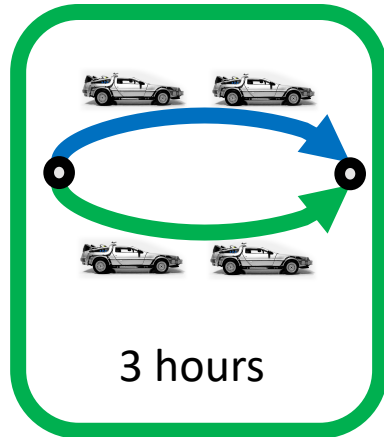Best option!
(Pareto optimal)

Selfish traffic
is like this!

Traffic

Total
Time          3.25 hours      3 hours      3.25 hours      4 hours

1 hour

1 hour

1 hour

1 hour

15 minutes per car

S

Travel Time

D

1 hour

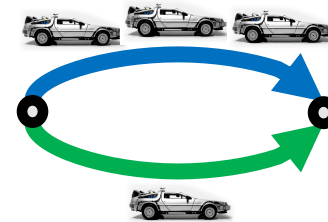Best option!
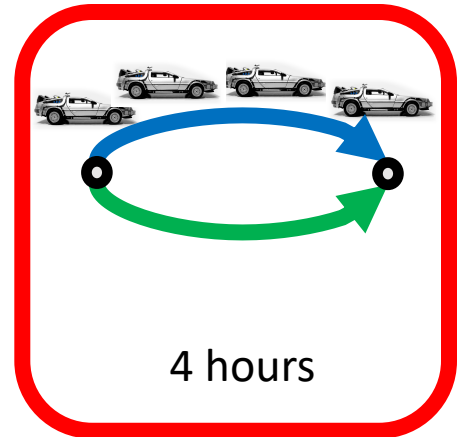(Pareto optimal)

Selfish traffic
is like this!

Traffic

Total
Time

3.25 hours

3 hours

3.25 hours

4 hours

1 hour

1 hour

1 hour

2 hours

15 minutes per car

Travel Time

1 hour

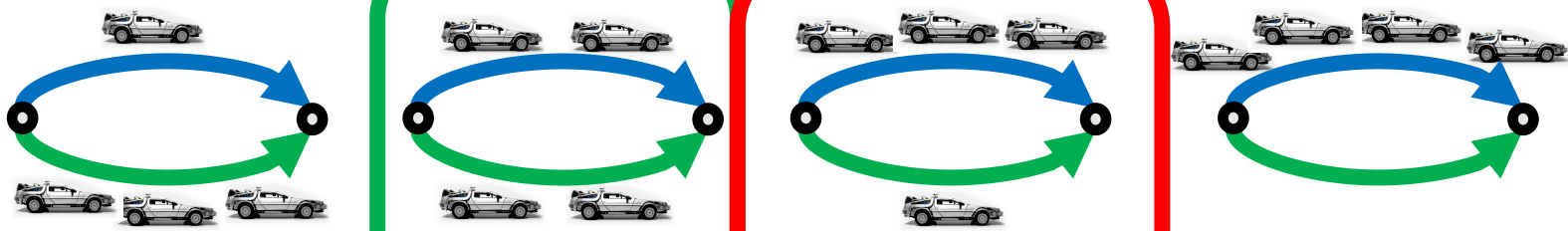S                                                    D

Best option!
(Pareto optimal)

Selfish traffic
is like this!

Traffic

Total
Time

3.25 hours        3 hours        3.25 hours        4 hours

45 minutes

45 minutes

90 minutes

1 hour

15 minutes per car

Travel Time

1 hour

S

D

Best option!
(Pareto optimal)

Traffic

Total
Time

3.25 hours

3 hours

3.25 hours

4 hours

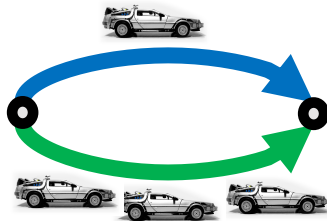30 minutes
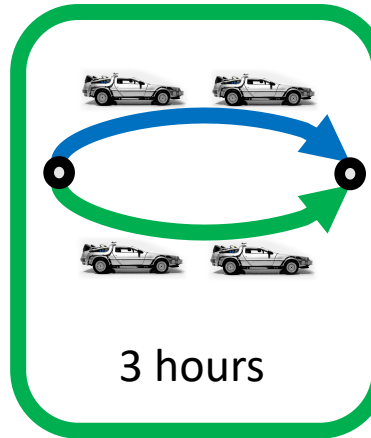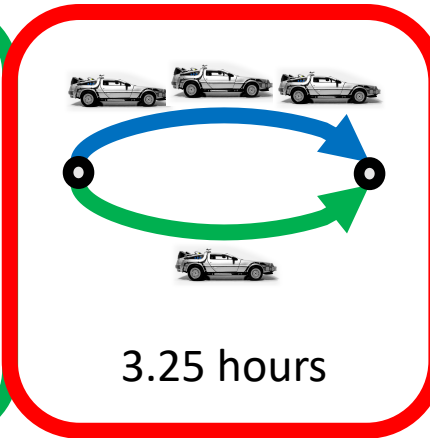
30 minutes

15 minutes per car

S

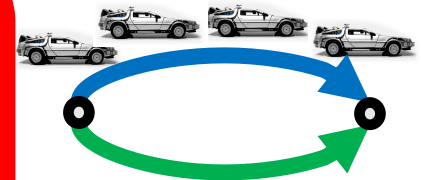Travel Time

D

1 hour

1 hour

1 hour

Best option!
(Pareto optimal)

Traffic

Total
Time

3.25 hours

3 hours

3.25 hours

4 hours

University of Colorado
Colorado Springs

University of Colorado
Boulder | Colorado Springs | Denver | Anschutz Medical Campus

Altruistic self-driving cars:

- Improve congestion
- Even if only some are altruistic
- Without making others worse off
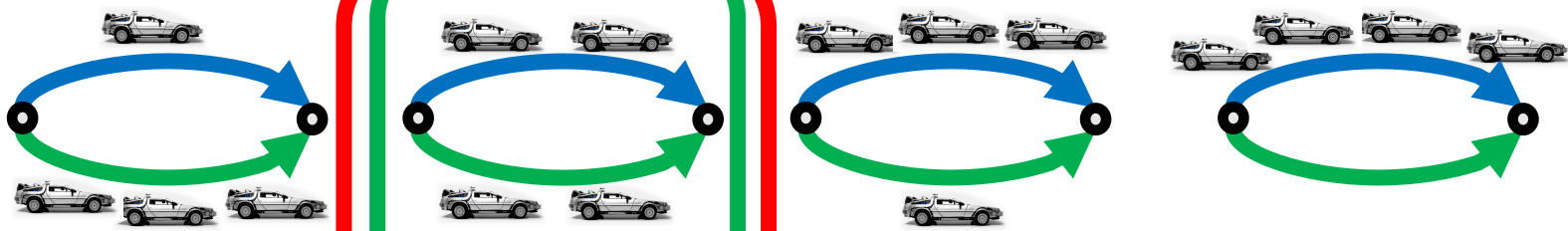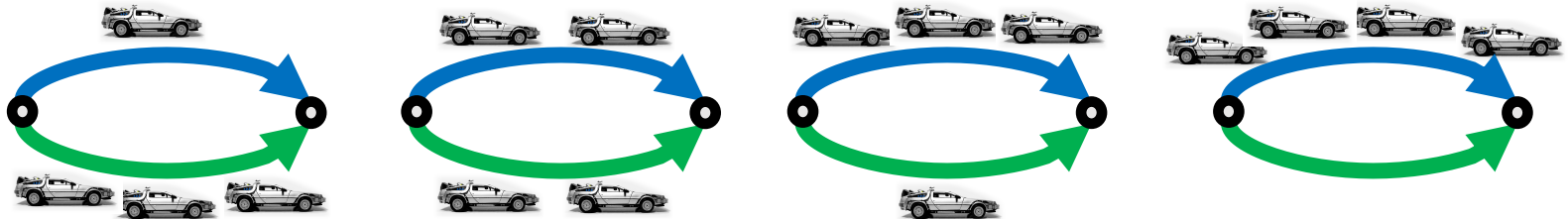- Unambiguously Ethical?



Traffic

Total
Time

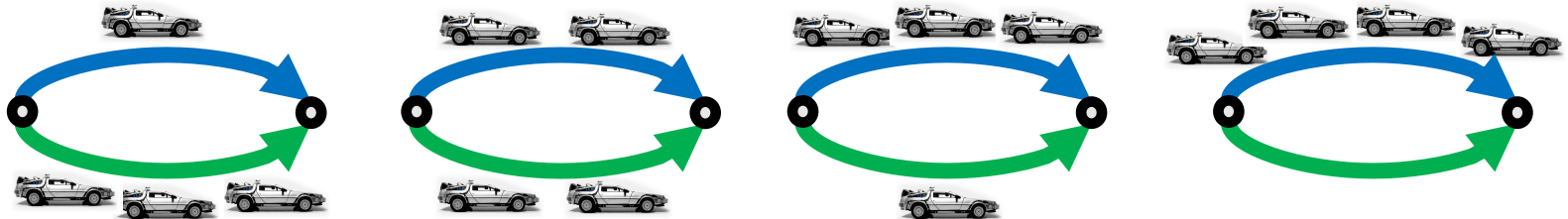3.25 hours          3 hours          3.25 hours          4 hours

Altruistic self-driving cars:

- Improve congestion
- Even if only some are altruistic
- Without making others worse off
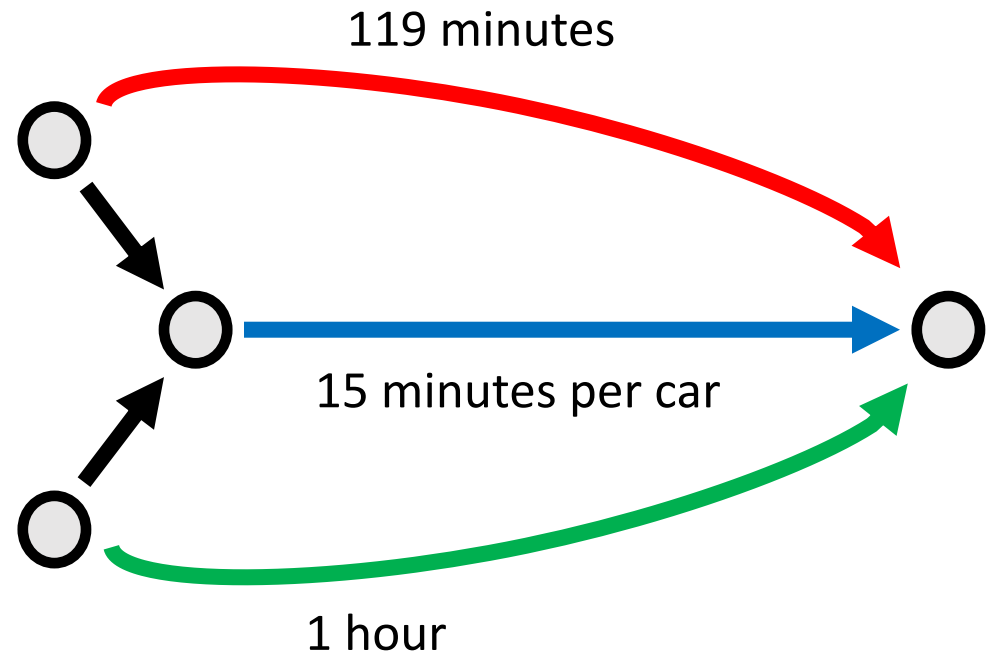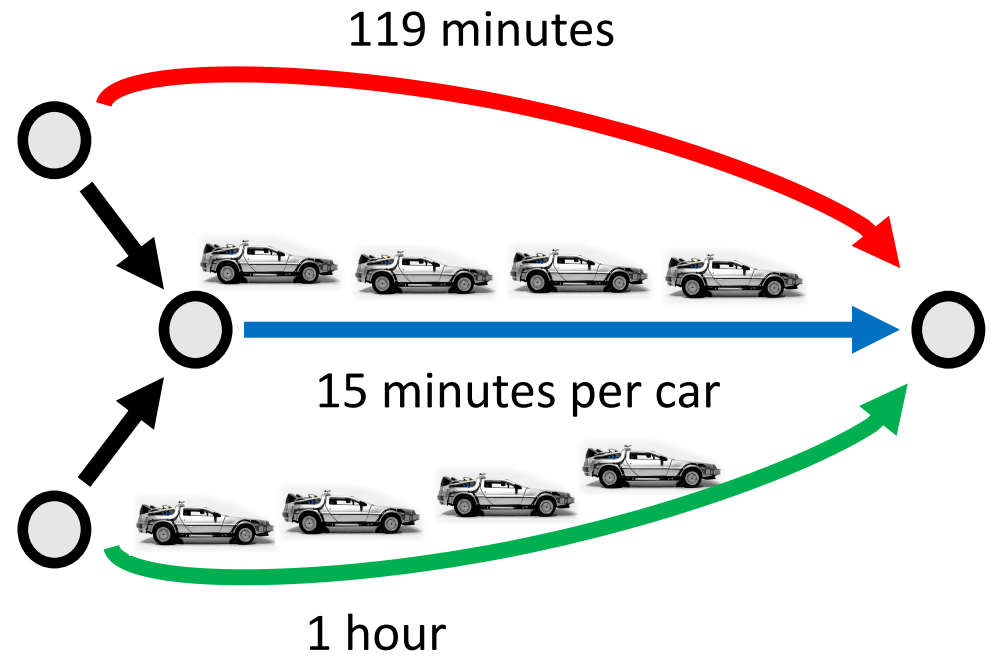- Unambiguously Ethical?
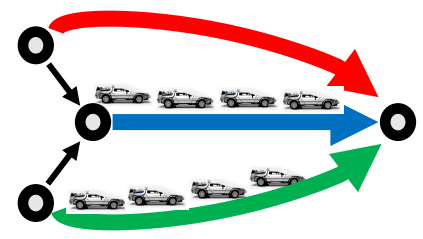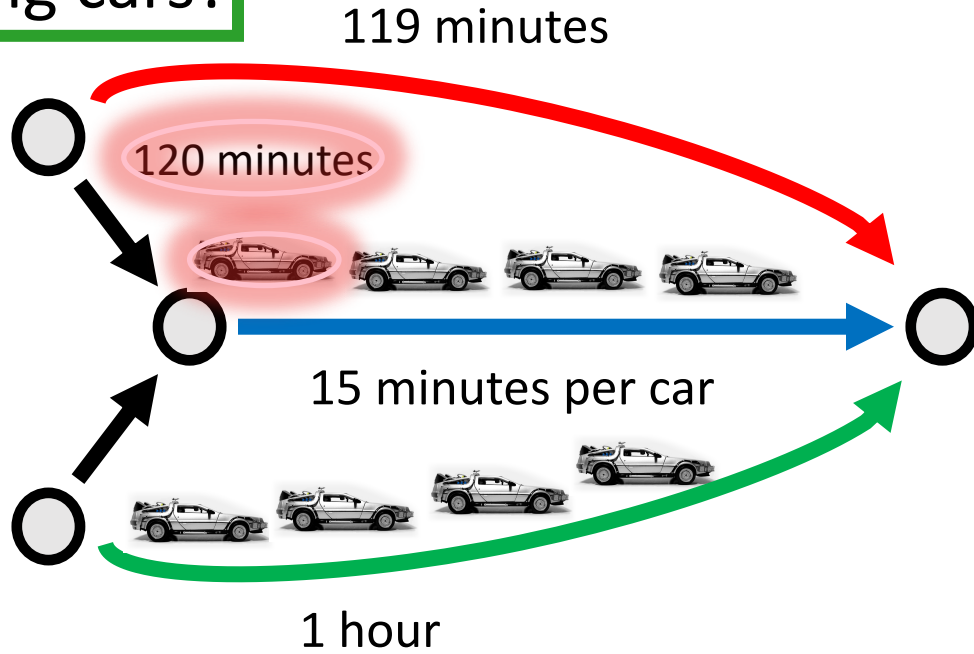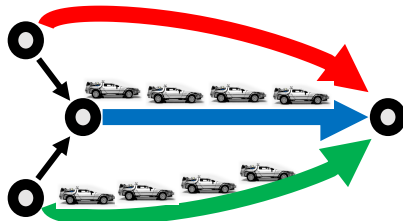
Traffic

Total
Time

| 3.25 hours | 3 hours | 3.25 hours | 4 hours |

119 minutes

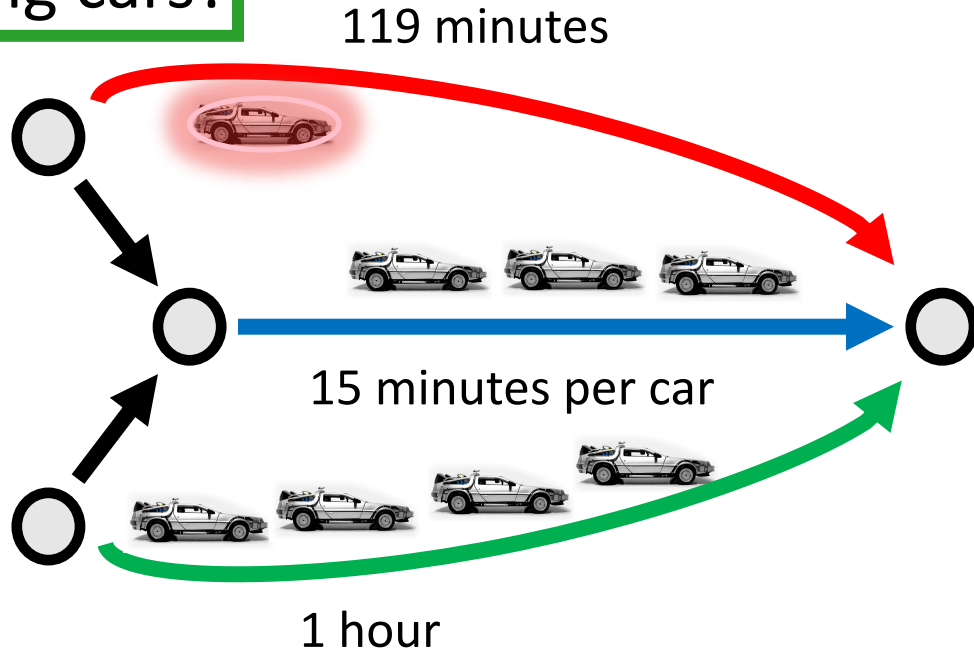15 minutes per car

1 hour

119 minutes

15 minutes per car

1 hour

Selfish Traffic

Total
Time

8 hours

# Altruistic self-driving cars?

119 minutes

120 minutes

15 minutes per car

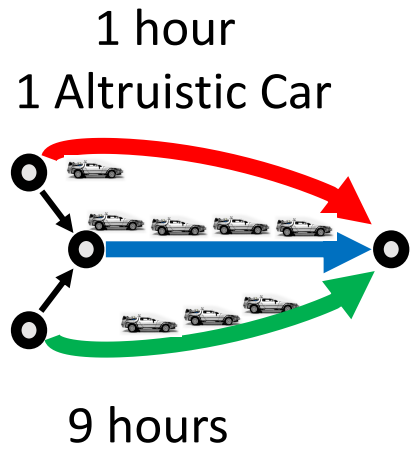1 hour
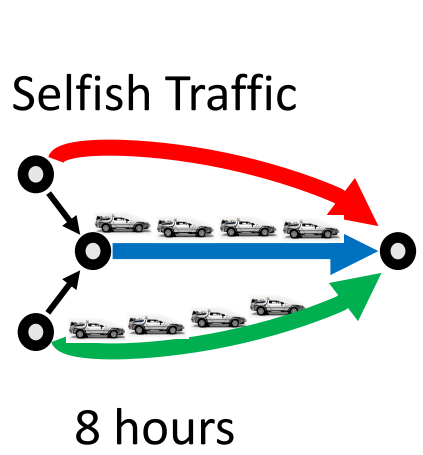
Selfish Traffic

Total
Time

8 hours

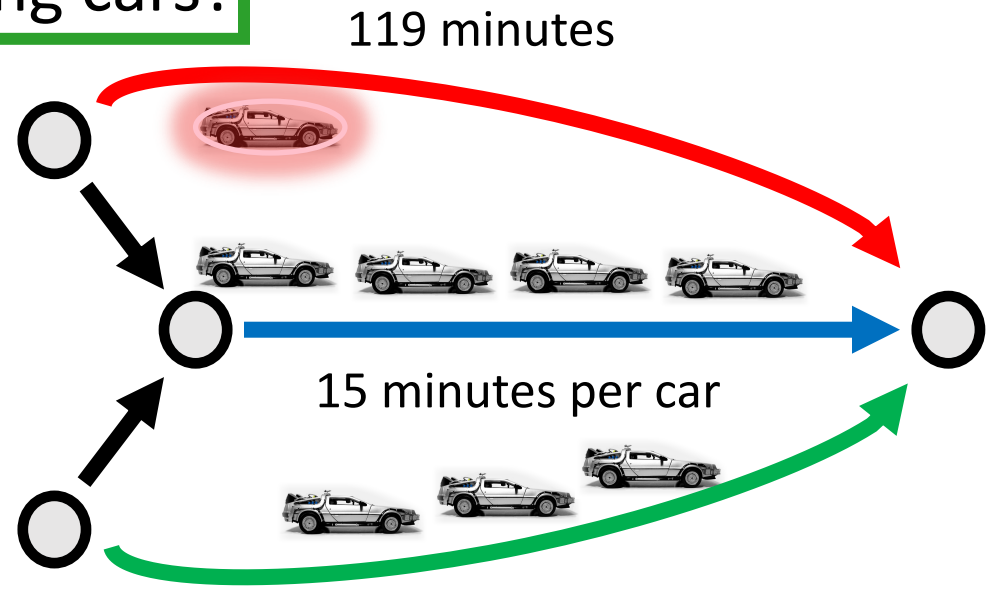# Altruistic self-driving cars?

119 minutes

15 minutes per car

1 hour

Selfish Traffic

Total
Time

8 hours

# Altruistic self-driving cars?

119 minutes

15 minutes per car

1 hour
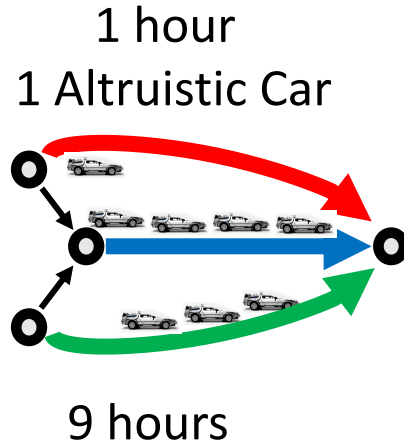1 Altruistic Car

Selfish Traffic

Total
Time

8 hours

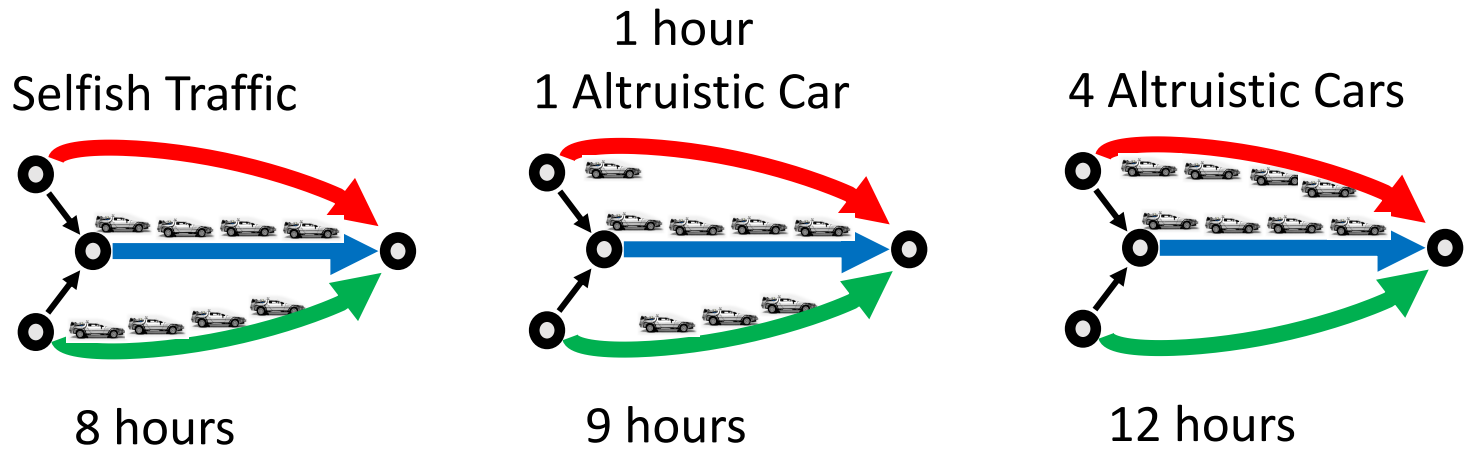9 hours

# Altruistic self-driving cars?

119 minutes

15 minutes per car

1 hour

Selfish Traffic

1 Altruistic Car

4 Altruistic Cars

Total
Time

8 hours

9 hours

12 hours

Decision design for socio-technical systems

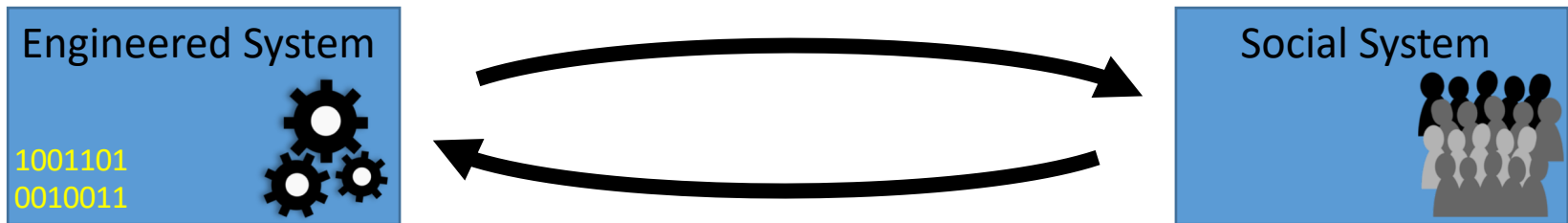|  | | |
|---|---|---|
| **Machine Policy** |  |  |
| Selfish | Not Optimal | Not Optimal |
| Altruistic | Inconsiderate? | Also potentially inefficient! |

If **Ethical** means **utilitarian**, then **sometimes altruism is good!**

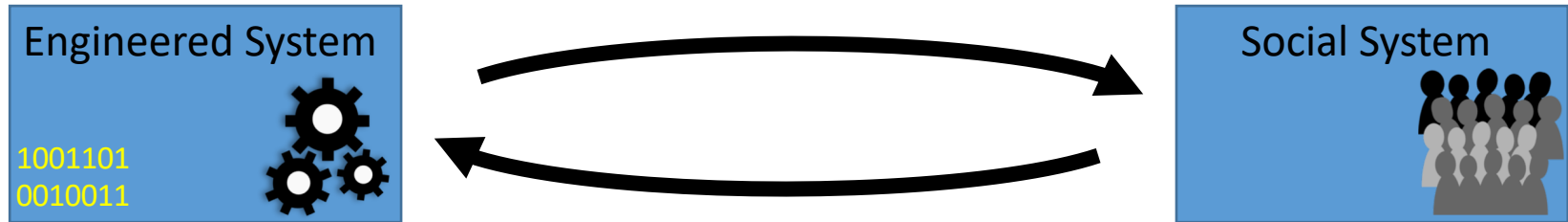If **Ethical** means **risk-averse**, then **machines should be selfish!**

Engineered System

1001101
0010011

Social System

**Needed Update:** Engineer machines to *interact* well

## What I'm doing

- Teaching CS4730/5730: Algorithmic Game Theory
- Designing CS4740: Social and Engineering Networks
    - Mathematics of *interaction*
- Public Lectures on responsible technological interaction

## What the community can do

- Engineering curricula need integrated social science
- Assume that interaction → counterintuitive outcomes