# Predicting Solar Irradiance Data Using Machine Learning

Matt Franks, Associate Principal
2019 Radiance Workshop
August 22, 2019
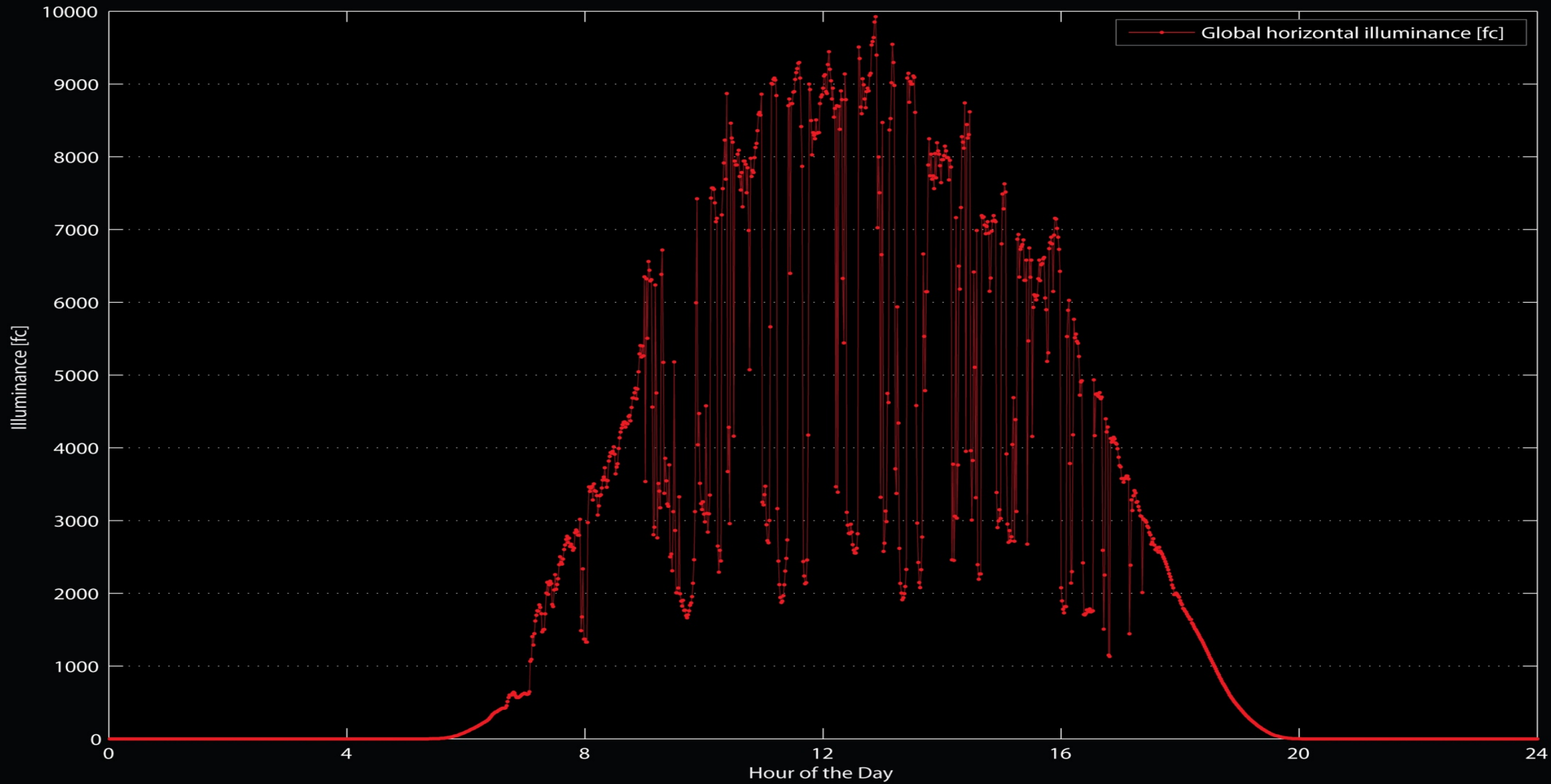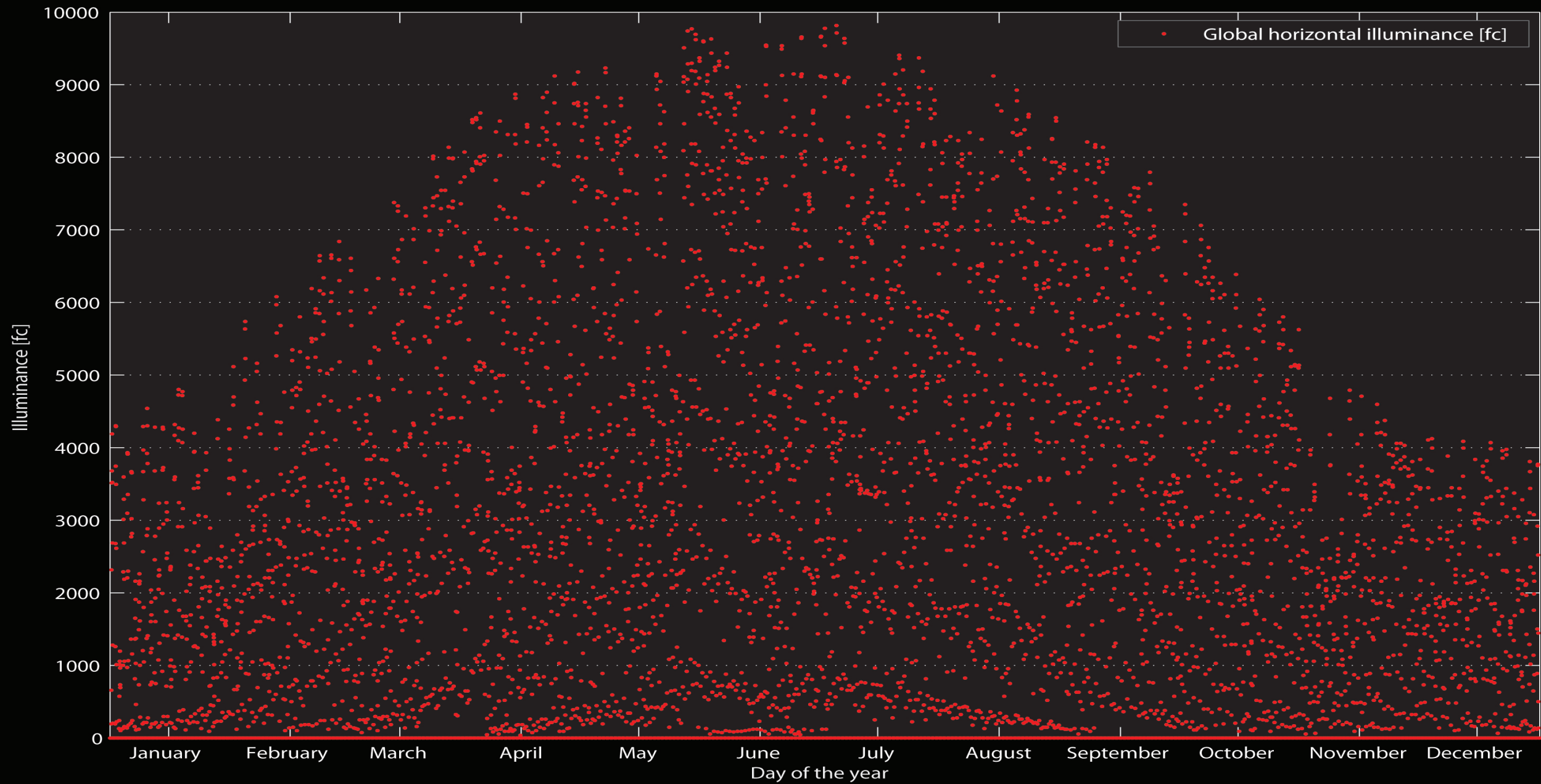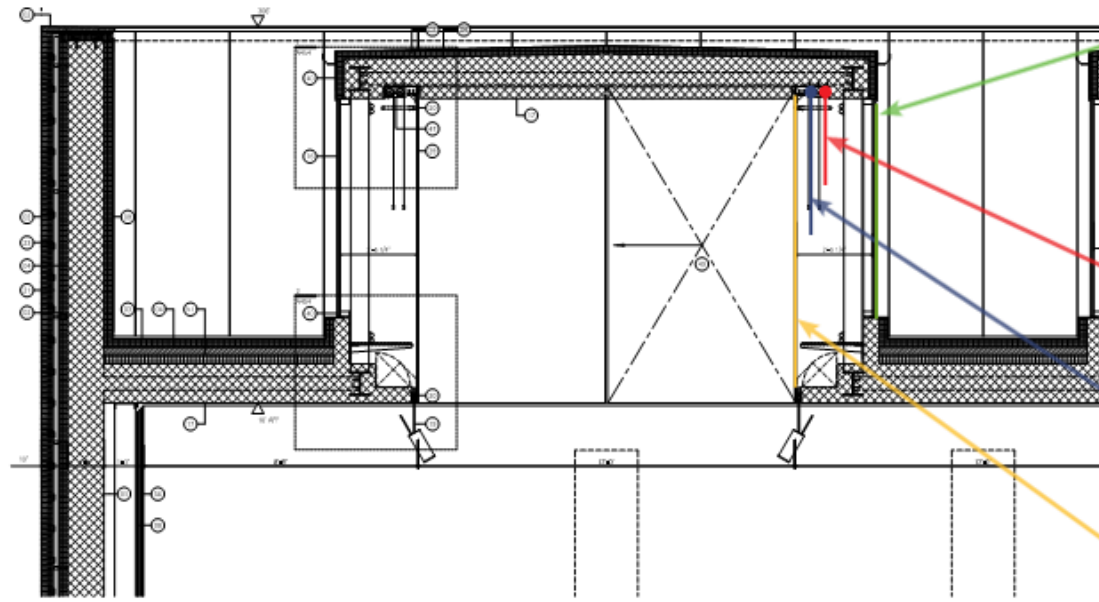
ARUP

One Day

One Year

Section Detail - Gallery 2



Section - Gallery 2



Keyplan

## Roof Monitor Layers:
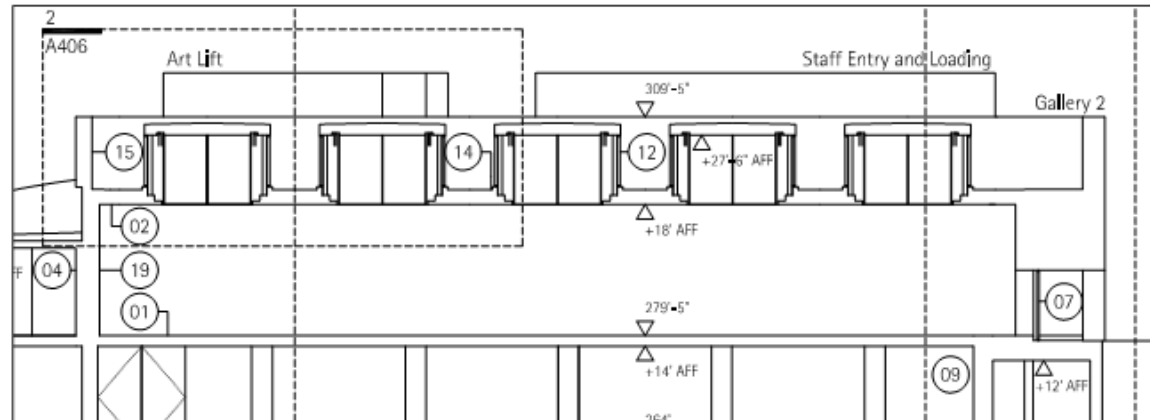
Low iron insulated glazing unit with laminated diffusing inner lite to provide ultra violet filter. Low iron glass is used to maximise colour rendering. Low-e coatings are as neutral in color terms as possible to maintain color rendering of the skylight glass unit of 97 or above. The laminated inner layer will be diffusing to mitigate direct sunlight penetration.
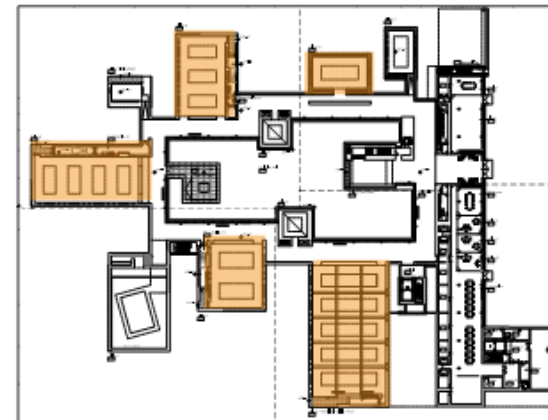
Motorised blackout roller shade to reduce daylight exposure outside museum open hours and allow for flexibility in the allowance of daylight into the gallery.The blackout shade should be provided with side-channels to eliminate light spill around the edges of the shade.

Motorised dimout roller shade to allow for reduction of light levels passing through the skylight system. The shade shall be an open-weave materials with 5% openness and a 10% to 15% visible light transmission, to be determined.

Interior diffusing glass to further diffuse directionality of light and obscure view of structure, roller shades and lighting. This shall be laminated with a diffusing interlayer, and be operable to allow easy access for maintenance. The interior glazing will have an acid-etch finish to reduce interior specular reflections.

## 3.6    Galleries 2, 4, 6, 8, and 10

### 3.6.1  Approach

Galleries 2, 4, 6, 8, and 10 will have similar daylight system designs, consisting of a roof monitor system. Refer to the architectural plans for the arrangement and dimensions of roof monitors in each gallery.

It is expected that Gallery 2 will generally be used to display parts of the permanent collection – typically a mixture of oil paintings, photographs, sculpture.

It is also recognized that Gallery 2 would at times be used for mixed media collections, which means that some works on paper may be displayed along with oil paintings and sculpture. Blackout, if required, is proposed to be provided by the deployment of roller shades in the roof monitors.
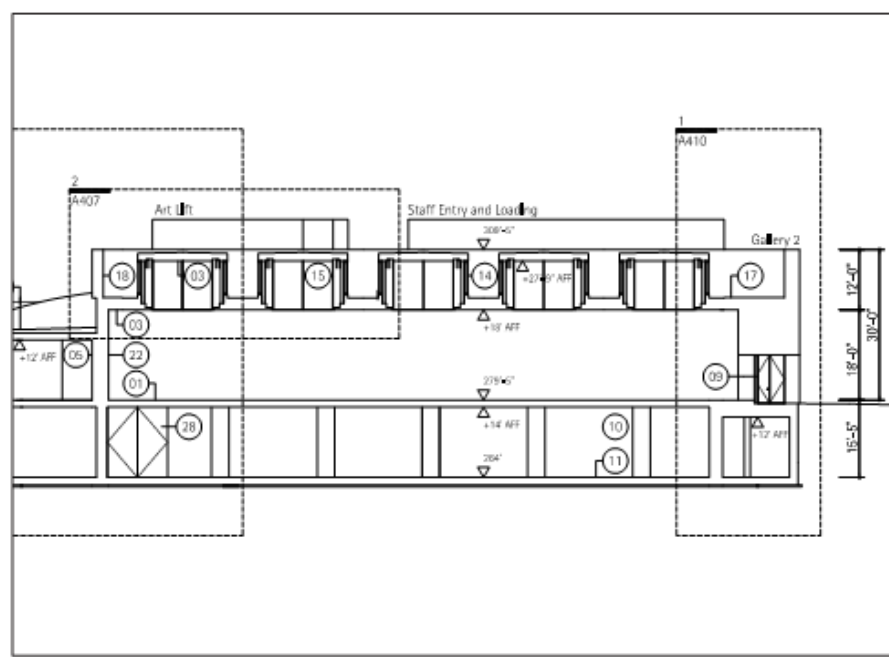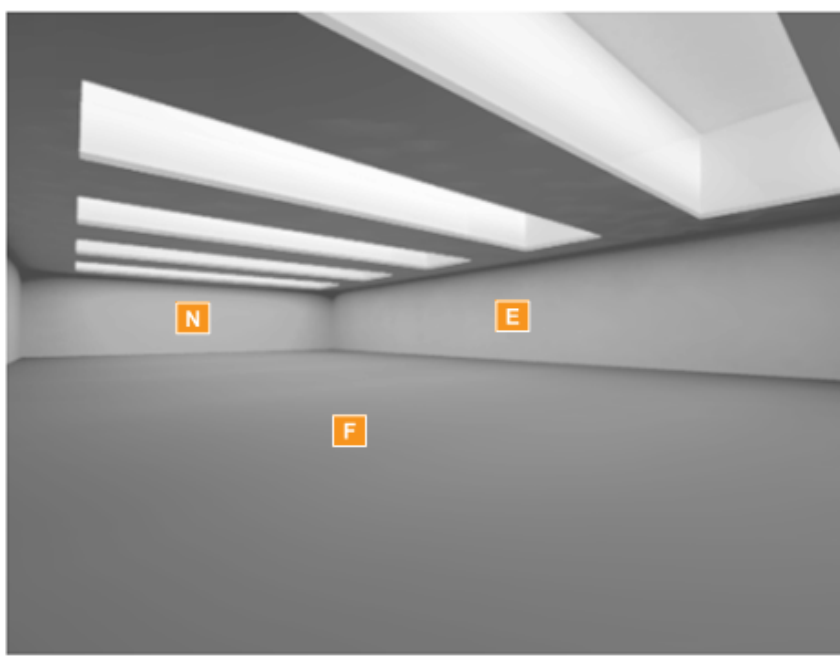
Galleries 4, 6, 8, and 10 will be used for more permanent exhibitions. It is understood that upon completion of construction Galleries 4 and 6 will exclude daylight due to the nature of their exhibits, however provisions for daylighting will be included in the design.

### 3.6.2  Proposed daylight system

The ceiling consists of a roof monitor system which introduces generous but controlled daylight into the gallery below. The images to the left illustrate the proposed system, which consists of a number of layers:

- Exterior vertical diffusing glass, running in the east-west direction
- Interior motorized blackout shade
- Interior motorized roller shade
- Interior diffusing Glass

These sets of layers will occur on both the north and south sides of the roof monitor. By allowing sunlight to be diffused through the layers of the southern glazing, and northern skylight to be transmitted and diffused through the north-facing glazing, the lighting conditions in the gallery will vary through out the day as sun position and weather patterns change.
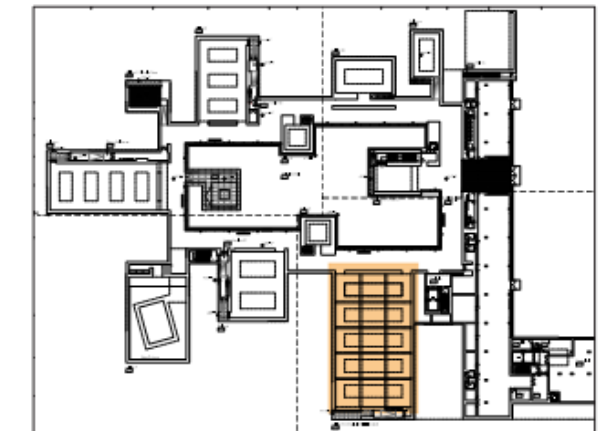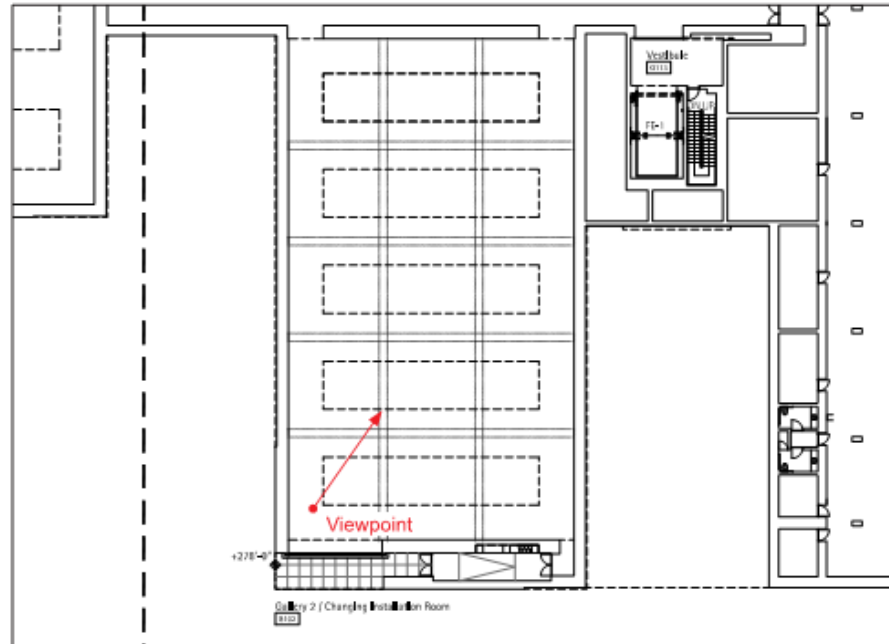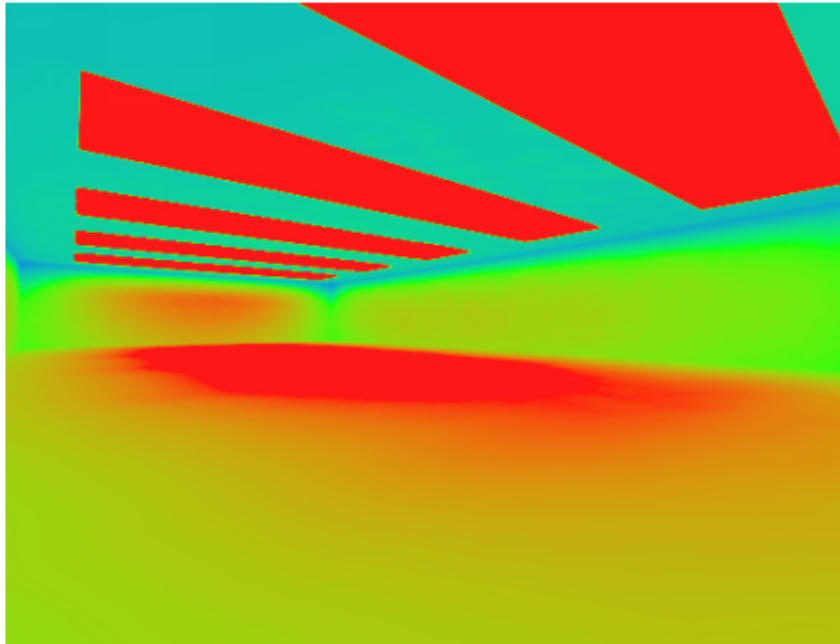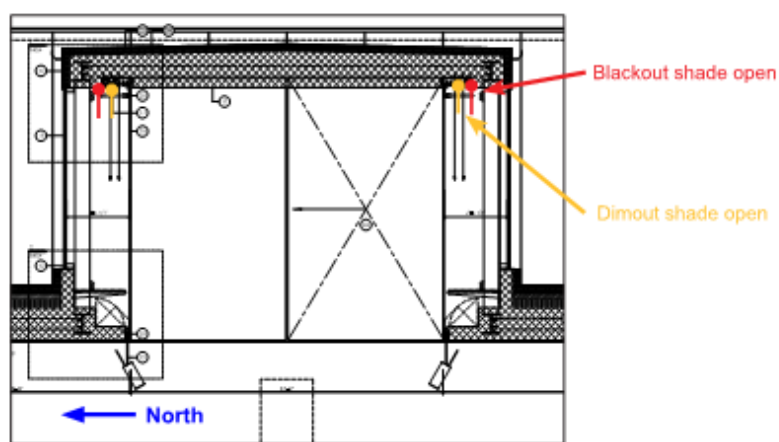
Viewpoint

## 4.4 Gallery 2

- Outer Glazing Transmittance: 53%
- Inner Glazing Transmittance: 64%
- Wall Reflectance: 75%
- Concrete Reflectance: 70%
- Floor Reflectance: 50%
- Calculation Time: 12:00 p.m. on date indicated
- Measurement points are at location indicated in images. Unless otherwise indicated images shown are for March 31, overcast conditions.

| Day | Weather | N (fc) | E (fc) | S (fc) | W (fc) | F (fc) |
|---|---|---|---|---|---|---|
| Mar 21 | Overcast | 114 | 93 | 106 | 94 | 143 |
| | Sunny | 351 | 257 | 220 | 262 | 388 |
| Jun 21 | Overcast | 119 | 100 | 116 | 101 | 153 |
| | Sunny | 266 | 200 | 215 | 203 | 298 |
| Dec 21 | Overcast | 72 | 59 | 65 | 59 | 91 |
| | Sunny | 149 | 119 | 92 | 118 | 186 |

*N, E, S, W, measurement points are on North, East, South, and West walls respectively; F measurement point is horizontal illuminance on floor. S and W points are not pictured.*


Keyplan

Gallery 2 - Roof Monitor Section

| Day | Weather | North (fc) | East (fc) | South (fc) | West (fc) | Floor (fc) |
|---|---|---|---|---|---|---|
| Mar 21, 12:00 p.m. | Overcast | 114 | 93 | 106 | 94 | 143 |
| | Sunny | 351 | 257 | 220 | 262 | 388 |

| | North (k-fc-hr) | East (k-fc-hr) | South (k-fc-hr) | West (k-fc-hr) | Floor (k-fc-hr) |
|---|---|---|---|---|---|
| Annual Cumulative Exposure | 457 | 350 | 362 | 353 | 524 |



Gallery 2 - Plan



Annual Illuminance Profile on East Wall at Point E

Gallery 2 - Roof Monitor Section

| Day | Weather | North (fc) | East (fc) | South (fc) | West (fc) | Floor (fc) |
|---|---|---|---|---|---|---|
| Mar 21, 12:00 p.m. | Overcast | 114 | 93 | 106 | 94 | 143 |
| | Sunny | 351 | 257 | 220 | 262 | 388 |

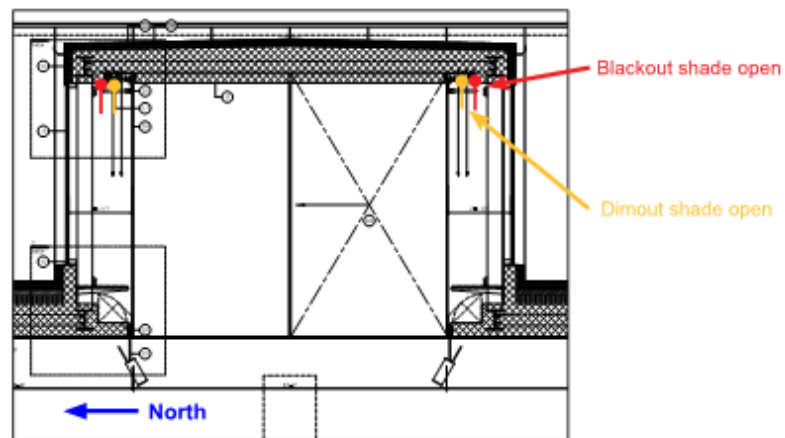| | North (k-fc-hr) | East (k-fc-hr) | South (k-fc-hr) | West (k-fc-hr) | Floor (k-fc-hr) |
|---|---|---|---|---|---|
| Annual Cumulative Exposure | 199 | 152 | 156 | 154 | 227 |



Gallery 2 - Plan



Annual Illuminance Profile on East Wall at Point E

Gallery 2 - Roof Monitor Section

| Day | Weather | North (fc) | East (fc) | South (fc) | West (fc) | Floor (fc) |
|---|---|---|---|---|---|---|
| Mar 21, 12:00 p.m. | Overcast | 33 | 27 | 31 | 28 | 42 |
| | Sunny | 123 | 79 | 66 | 79 | 116 |

| | North (k-fc-hr) | East (k-fc-hr) | South (k-fc-hr) | West (k-fc-hr) | Floor (k-fc-hr) |
|---|---|---|---|---|---|
| Annual Cumulative Exposure | 59 | 46 | 47 | 46 | 68 |



Gallery 2 - Plan



Annual Illuminance Profile on East Wall at Point E

Gallery 2 - Roof Monitor Section

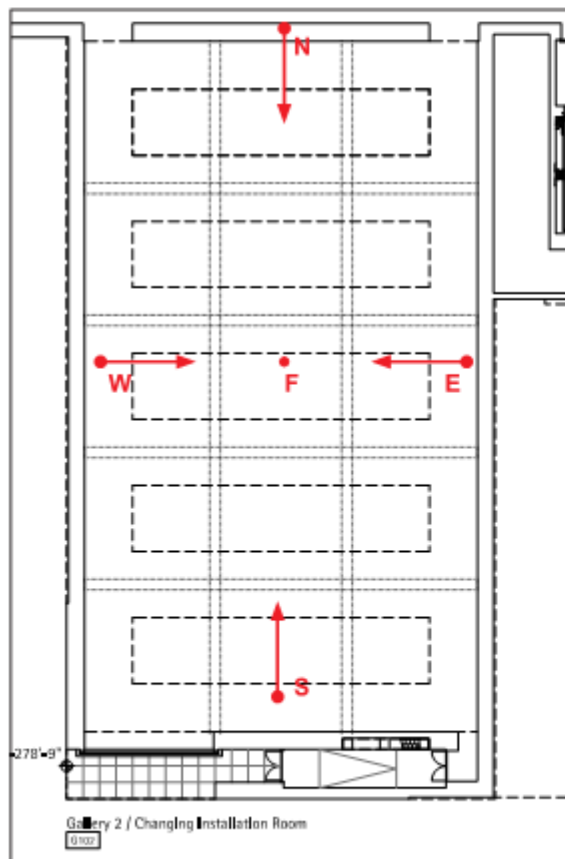| Day | Weather | North (fc) | East (fc) | South (fc) | West (fc) | Floor (fc) |
|---|---|---|---|---|---|---|
| Mar 21, 12:00 p.m. | Overcast | 90 | 60 | 50 | 61 | 91 |
| | Sunny | 345 | 212 | 138 | 214 | 293 |

| | North (k-fc-hr) | East (k-fc-hr) | South (k-fc-hr) | West (k-fc-hr) | Floor (k-fc-hr) |
|---|---|---|---|---|---|
| Annual Cumulative Exposure | 169 | 108 | 81 | 109 | 163 |



Gallery 2 - Plan



Annual Illuminance Profile on East Wall at Point E

# Why might reality be different than what was predicted?

- Real reflectances differ from those assumed
- Dirt more or less than assumed
- Constructed dimensions differ from design
- Inaccuracy of calculation methods

# Why might reality be different than what was predicted?

- Real reflectances differ from those assumed     <span style="color:red">**+/- 5%**</span>
- Dirt more or less than assumed     <span style="color:red">**+/- 5%**</span>
- Constructed dimensions differ from design     <span style="color:red">**+/- 1%**</span>
- Inaccuracy of calculation methods     <span style="color:red">**+/- 5%**</span>

## 2.3 Illuminance Distribution

The images on this page show comparisons between the computer model daylight distribution simulation, both in greyscale and falsecolor luminance.

The luminance distribution images show reasonable uniformity as well as agreement with the computer simulated distribution. Note that the slight dropoff in the center of wall on the right side of the image is due to the model construction, which consists of a mirror to replicate the appearance of two additional clerestories.



Computer Simulation - Gallery 2 - Greyscale



Photo - Gallery 2 Model



Computer Simulation - Gallery 2 - Falsecolor



Photo - Gallery 2 Model - Falsecolor

Gallery 2 Illuminance - Calculated and Measured

Legend:
- Gallery 2 Calculated
- Gallery 2 Measured

Chart regions:
- No Shades Installed
- Shade on South Clerestory
- Shades on South and North Clerestories

Y-axis: Illuminance / footcandles (0.0 to 400.0)
X-axis: Jan through Dec

## 2 Gallery 2

### 2.1 Illuminance Levels - Annual

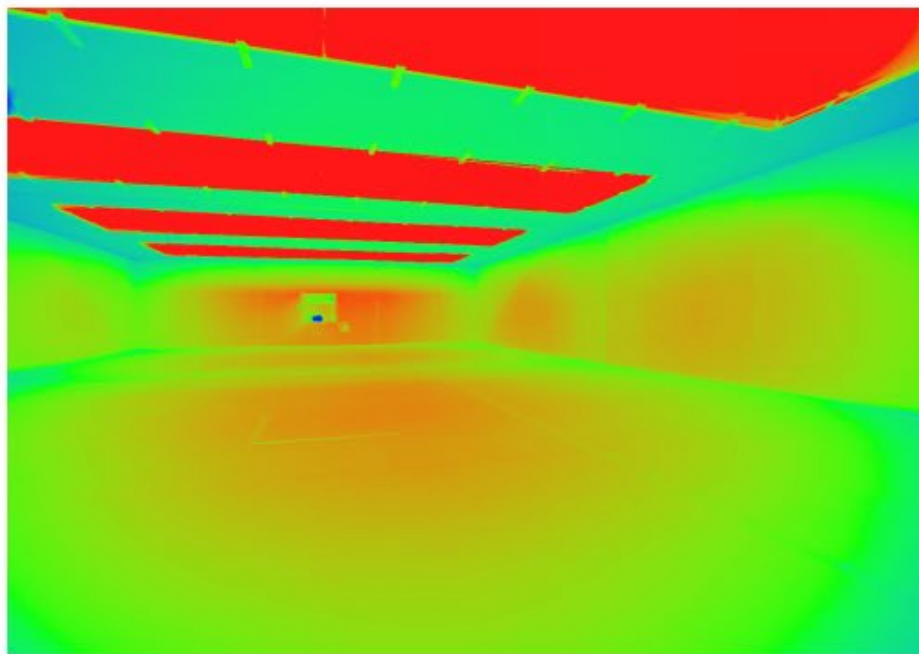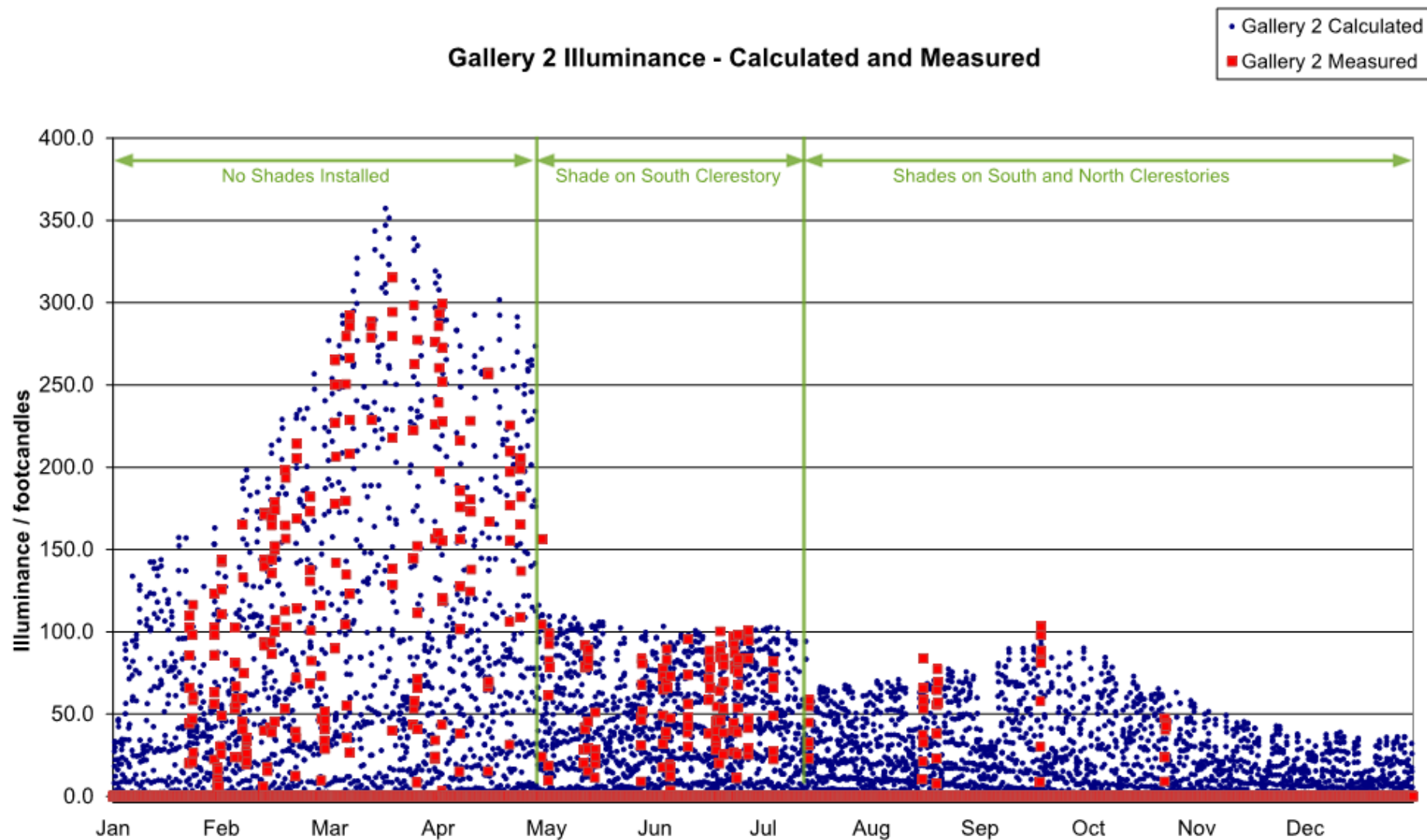The scatter plot on the left side of this pages shows a comparison between the hourly illuminance data calculated for the north wall of Gallery 2 for each hour of the year based on typical weather data from the November 15, 2012 daylighting report (blue dots), overlayed with hourly data measured from the Gallery 2 model on days that measurements was possible (red squares).

Indicated on the scatter plot is the times that the three different shade configurations were installed in the model.

- No shades from January to April 30.
- Shades on only the south clerestory from May 1 to July 14.
- Shades on both the north and south clerestories for the remainder of the year.

The shade material used in the model was the shade material currently specified, Mermet Screen Vision:

- 10% openness
- 29% visible light transmittance
- white color

The general trend of the data indicates fairly close correlation between the computer model and the measured illuminance, with illuminance peaks at similar levels.

Gallery 9 Illuminance - Calculated and Measured

- Gallery 9 Calculated
- Gallery 9 Measured

# 3 Gallery 9

## 3.1 Illuminance Levels - Annual

The scatter plot on the left side of this pages shows a comparison between the hourly illuminance data calculated for the west wall of Gallery 9 for each hour of the year based on typical weather data from the November 15, 2012 daylighting report (blue dots), overlayed with hourly data measured from the Gallery 9 model on days that measurements was possible (red squares).

The general trend of the data indicates fairly close correlation between the computer model and measurements in relative terms, however it can be seen from scatter plot that the gallery 9 measurements are in the range of 30% lower than predicted.

There are several factors that may be contributing to the difference between the calculated and measured values. These are discussed on the following pages.

# Communicating the Qualitative and Quantitative in Museum Daylighting

Kristen N. Garibaldi

**2017** INTERNATIONAL RADIANCE WORKSHOP

PORTLAND, OREGON
AUGUST 23, 2017

**ARUP**

# Problem (?)

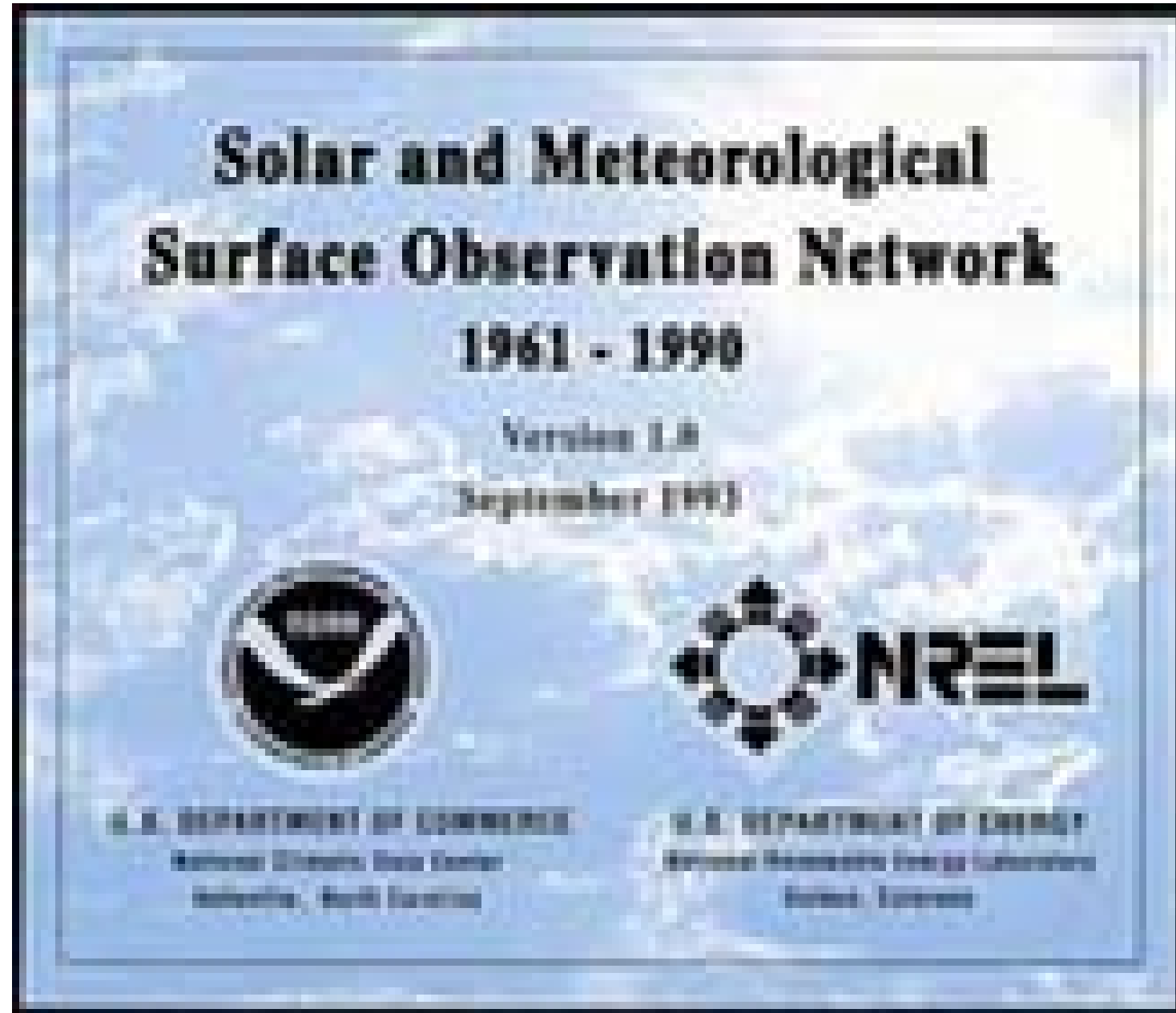- It is difficult to measure direct and diffuse illuminance (irradiance) separately.
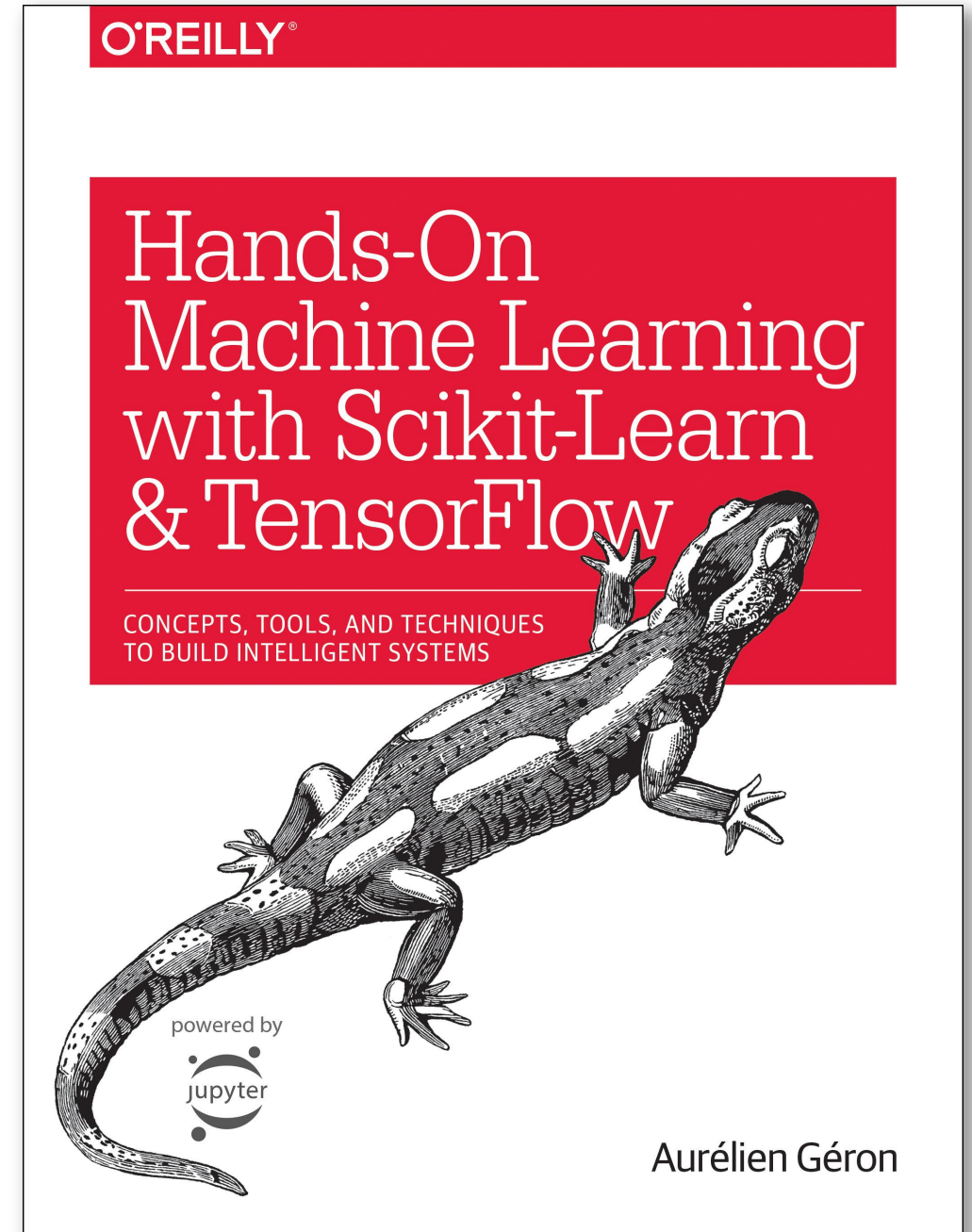
$200
1" diameter
2 ounces

$6,000
6" diameter
2 lbs

But…



Solar and Meteorological
Surface Observation Network
1961 - 1990

Version 1.0

September 1993

| DATA SET | |
|---|---|
| WHAT | Weather Data |
| FORMAT | TMY |
| LOCATION | Local Storage |

## 236 Folders

My Computer – N-YLTCND7174M1Y > Desktop > weather files > SAMSON (US)

| Name | Date modified | Type |
|---|---|---|
| 24284 | 11/16/2017 5:05 PM | File folder |
| 25308 | 11/16/2017 5:05 PM | File folder |
| 25339 | 11/16/2017 5:05 PM | File folder |
| 25501 | 11/16/2017 5:05 PM | File folder |
| 25503 | 11/16/2017 5:05 PM | File folder |
| 25624 | 11/16/2017 5:05 PM | File folder |
| 25713 | 11/16/2017 5:05 PM | File folder |
| 26411 | 11/16/2017 5:05 PM | File folder |
| 26415 | 11/16/2017 5:05 PM | File folder |
| 26425 | 11/16/2017 5:05 PM | File folder |
| 26451 | 11/16/2017 5:05 PM | File folder |
| 26510 | 11/16/2017 5:05 PM | File folder |
| 26528 | 11/16/2017 5:05 PM | File folder |
| 26533 | 11/16/2017 5:05 PM | File folder |
| 26615 | 11/16/2017 5:05 PM | File folder |
| 26616 | 11/16/2017 5:05 PM | File folder |

## 30 Files per Folder

| Clipboard | Organize | New |
|---|---|---|

My Computer – N-YLTCND7174M1Y > Desktop > weather files > SAMSON (US) > 93193

| Name | Date modified | Type |
|---|---|---|
| 93193_61.Z | 7/28/1993 3:38 AM | z Archive |
| 93193_62.Z | 7/28/1993 3:38 AM | z Archive |
| 93193_63.Z | 7/28/1993 3:39 AM | z Archive |
| 93193_64.Z | 7/28/1993 3:39 AM | z Archive |
| 93193_65.Z | 7/28/1993 3:39 AM | z Archive |
| 93193_66.Z | 7/28/1993 3:40 AM | z Archive |
| 93193_67.Z | 7/28/1993 3:40 AM | z Archive |
| 93193_68.Z | 7/28/1993 3:41 AM | z Archive |
| 93193_69.Z | 7/28/1993 3:41 AM | z Archive |
| 93193_70.Z | 7/28/1993 3:41 AM | z Archive |
| 93193_71.Z | 7/28/1993 3:42 AM | z Archive |
| 93193_72.Z | 7/28/1993 3:42 AM | z Archive |
| 93193_73.Z | 7/28/1993 3:42 AM | z Archive |
| 93193_74.Z | 7/28/1993 3:43 AM | z Archive |

## 1 header row (same for 30 files)

26451_67

```
1   26451 ANCHORAGE              AK -10  N61 10  W150 01     35
2   67 1 1 1    0    0    0 ?0    0 ?0    0 ?0  9  9  -8.3  -9.4 92 996  10 1.5 24.1   910 0999999999   699999.  41 0
3   67 1 1 2    0    0    0 ?0    0 ?0    0 ?0 10 10  -6.7  -8.9 84 996   0  .0 24.1   760 0999999999   699999.  41 0
4   67 1 1 3    0    0    0 ?0    0 ?0    0 ?0 10 10  -7.2  -8.9 88 996   0  .0 24.1   700 0999999999   699999.  41 0
5   67 1 1 4    0    0    0 ?0    0 ?0    0 ?0 10 10  -6.7  -8.9 84 996   0  .0 24.1   700 0999999999   699999.  41 0
6   67 1 1 5    0    0    0 ?0    0 ?0    0 ?0  8  7  -6.7  -8.3 88 995  30 2.1 24.1  2740 0999999999   699999.  41 0
7   67 1 1 6    0    0    0 ?0    0 ?0    0 ?0  7  6  -8.9 -10.0 92 995   0  .0 24.1  2740 0999999999   699999.  41 0
8   67 1 1 7    0    0    0 ?0    0 ?0    0 ?0  9  8 -10.0 -12.8 80 995   0  .0 24.1  1220 0999999999   699999.  41 0
9   67 1 1 8    0    0    0 ?0    0 ?0    0 ?0 10 10  -8.9 -11.1 84 995 360 1.5 24.1  1220 0999999999   699999.  41 0
10  67 1 1 9    0    0    0 ?0    0 ?0    0 ?0 10 10  -8.9 -10.6 88 996   0  .0 72.4  1520 0999999999   699999.  41 0
11  67 1 1 10  37  802    0 G5    1 G4    0 G5 10 10  -4.4  -6.7 85 996 280 2.1 80.5  1680 0999999999    6 .059  41 0
12  67 1 1 11  96 1415   25 G5    5 G4   25 G5 10 10  -6.1  -7.8 88 996 280 1.5 24.1  1680 0999999999    6 .059  41 0
13  67 1 1 12 139 1415   30 G5    5 G4   29 G5 10 10  -6.1  -7.2 92 996 100 2.1 32.2  1680 0999999999    6 .059  41 0
14  67 1 1 13 142 1415   56 G5    5 G4   56 G5 10 10  -5.6  -7.2 88 995 110  .5 12.9  1220 0999999999    6 .059  41 0
15  67 1 1 14 102 1415   37 G5    2 G4   36 G5 10 10  -5.6  -7.2 88 995   0  .0  8.0   180 0999999999    6 .059  41 0
16  67 1 1 15  42  943    6 G5    0 G4    6 G5 10 10  -5.6  -7.2 88 994 320 1.5  6.4   210 0999999999    6 .059  41 0
17  67 1 1 16   0    0    0 ?0    0 ?0    0 ?0 10 10  -5.6  -7.2 88 994 350 2.1  2.4   240 0999999999   699999.  41 0
18  67 1 1 17   0    0    0 ?0    0 ?0    0 ?0 10 10  -5.6  -7.2 88 993 340 1.0  3.2   430 0999999999   699999.  41 0
```

## 8760 rows of data, around 10 relavent columns

## 62 million total rows

# Machine Learning!

- "the science and art of programming computers so they can learn from data." (Geron)

- Machine learning uses data to "learn" and predict outcomes rather than using explicit algorithms or rules, and works well for problems that have no known algorithm based solution, but have lots of available data to learn from.
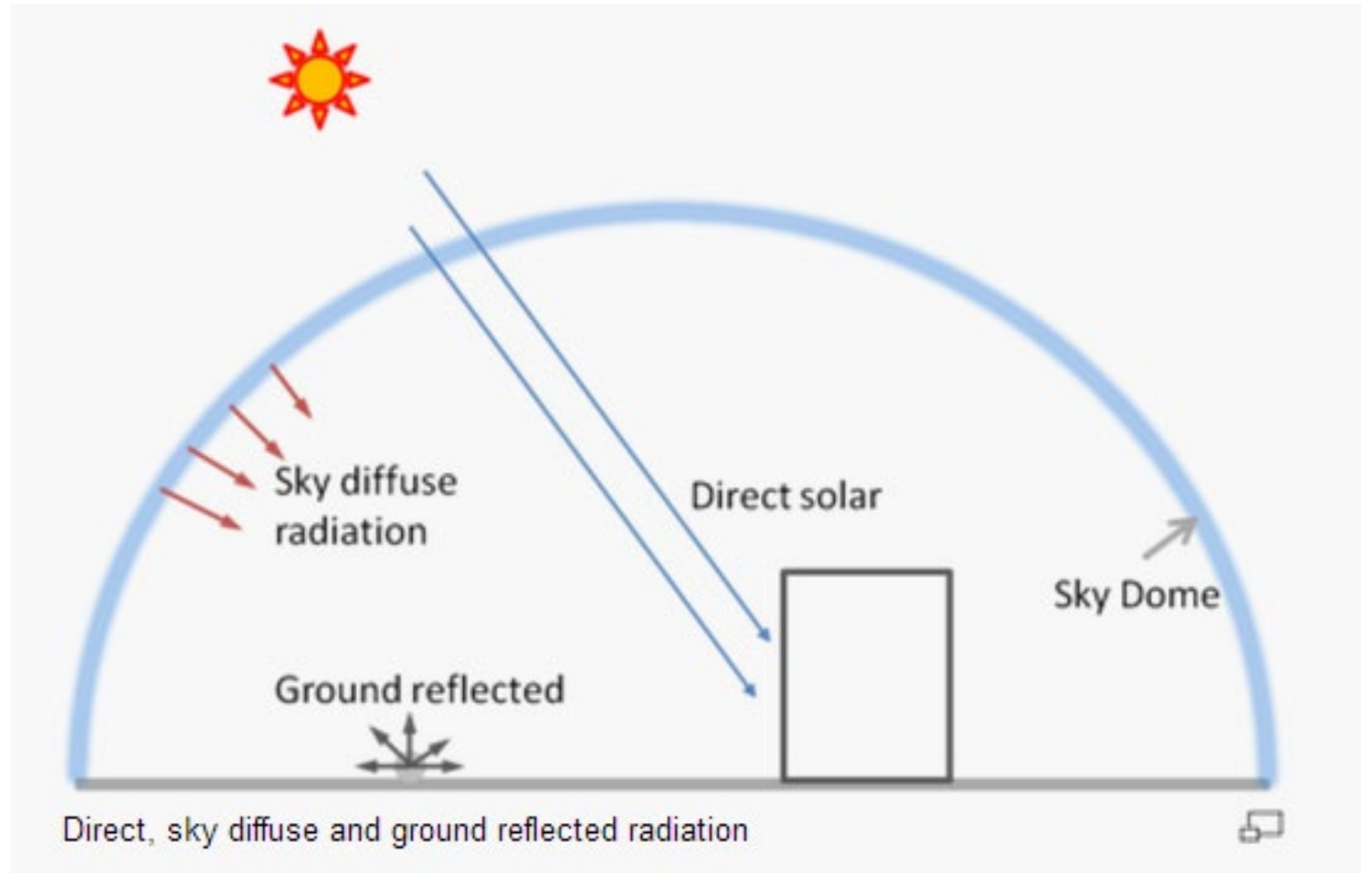
O'REILLY®

Hands-On Machine Learning with Scikit-Learn & TensorFlow

CONCEPTS, TOOLS, AND TECHNIQUES TO BUILD INTELLIGENT SYSTEMS

powered by
jupyter

Aurélien Géron

# When to use machine learning:

1. Tasks involve a function that maps well-defined inputs to well-defined outputs
2. Large (digital) datasets exist or can be created containing input-output pairs
3. Tasks provide clear feedback with clearly definable goals and metrics
4. No long chains of logic or reasoning that depend on diverse background knowledge or common sense
5. Tasks do not require detailed explanations for how the decision was made
6. Tasks have a tolerance for error and no need for provably correct or optimal solutions
7. The phenomenon or function being learned should not change rapidly over time
8. No specialized dexterity, physical skills, or mobility is required

From "What Can Machine Learning Do? Workforce Implications" Erik Brynjolfsson and Tom Mitchell, Science Magazine, Dec 22, 2017

# Hypothesis

- Data normally used:
  - Month
  - Day
  - Hour
  - Latitude
  - Longitude
  - Direct Illuminance (DIR)
  - Diffuse Illuminance (DIF)
- Data we also have:
  - Global Illuminance (GLOB)



Direct, sky diffuse and ground reflected radiation

$$DIR + DIF = GLOB$$

# Hypothesis

- Data normally used:
  - Month
  - Day
  - Hour
  - Latitude
  - Longitude
  - Direct Illuminance (DIR)
  - Diffuse Illuminance (DIF)
- Data we also have:
  - Global Illuminance (GLOB)

- If we have:
  - Month
  - Day
  - Hour
  - Latitude
  - Longitude
  - Global Illuminance (GLOB)
- Can we predict:
  - Direct Illuminance (DIR)
  - Diffuse Illuminance (DIF)

# Process

**this part only done once**

**Source Data**

Data Files:
1 header row, 8760 data rows
30 files per folder (location)
Files are in zip format
236 folders

→

**Parse Data**

1. unzip

2. split off header (add to data rows?)

→

**Store Data**

one big file? (~8.2gb, 62 million rows)

mysql database?

**this part only done for every project / location / measurement period**

**Select Data**

- specific location
- lat/long range
- ???

→

**Train ML Model**

- lat
- long
- sm
- date
- time
- global horiz
- direct normal
- diffuse horizontal

→

**ML Model**

→

**Hourly Results**

- direct normal
- diffuse horizontal

**Measured Data**

- lat
- long
- sm
- date
- time
- global horiz

# Tools

# Step 2 - Extracting data in a manner ready to use

The Python Data Analysis Library (pandas) has functions specifically to query and extract data from SQL databases in a data structure. The function below queries the first 12 rows where solar altitude is greater than zero and the state is 'NY'.

```python
In [24]: conn = sqlite3.connect("weather_database.db")
```

```python
In [29]: df = pd.read_sql_query("SELECT * FROM weatherdata WHERE solar_altitude > 0 AND state = 'NJ' LIMIT 12;", conn)
```
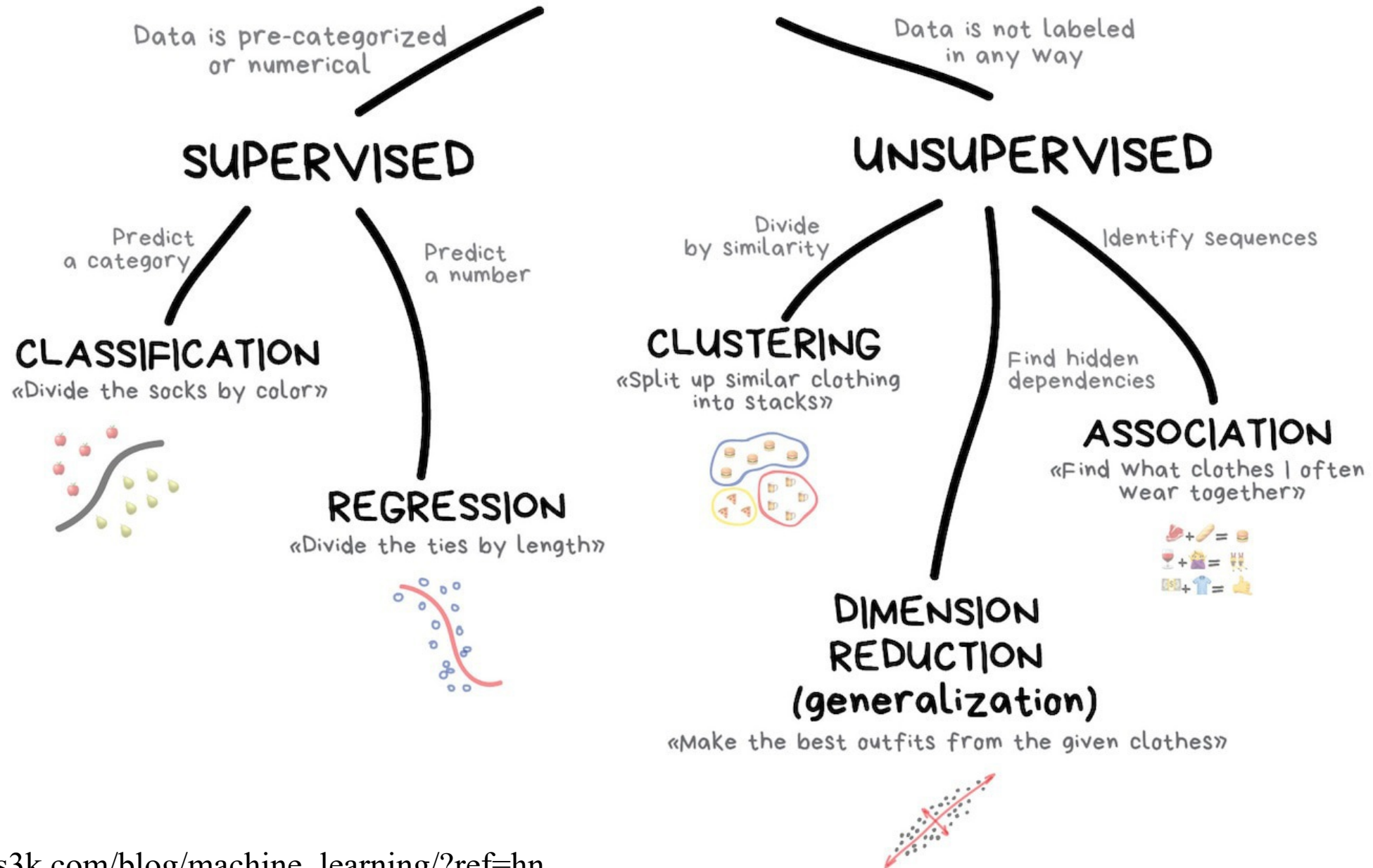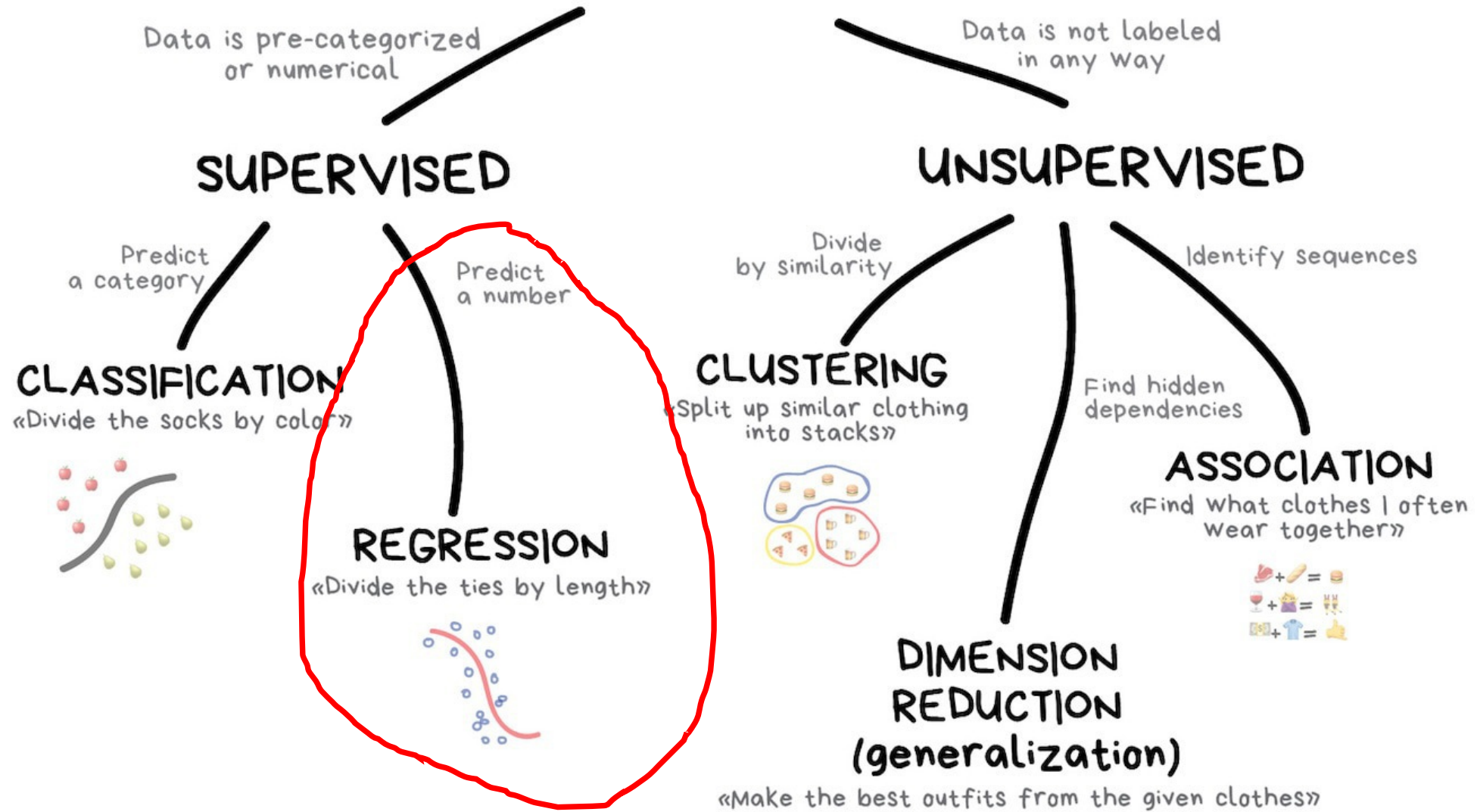
```python
In [30]: df
```

Out[30]:

| | station_id | city | state | timezone | lat | long | sm | elev | julian_day | yearhour | ... | month | day | hour | glob_horiz | dir_norm | dif_horiz | dir_h |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 8 | ... | 1 | 1 | 8 | 7.0 | 5.0 | 6.0 | |
| 1 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 9 | ... | 1 | 1 | 9 | 31.0 | 4.0 | 31.0 | |
| 2 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 10 | ... | 1 | 1 | 10 | 68.0 | 6.0 | 66.0 | |
| 3 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 11 | ... | 1 | 1 | 11 | 68.0 | 2.0 | 68.0 | |
| 4 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 12 | ... | 1 | 1 | 12 | 89.0 | 7.0 | 86.0 | |
| 5 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 13 | ... | 1 | 1 | 13 | 120.0 | 7.0 | 117.0 | |
| 6 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 14 | ... | 1 | 1 | 14 | 83.0 | 5.0 | 81.0 | |
| 7 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 15 | ... | 1 | 1 | 15 | 107.0 | 1.0 | 106.0 | |
| 8 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 16 | ... | 1 | 1 | 16 | 53.0 | 1.0 | 53.0 | |
| 9 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 2 | 32 | ... | 1 | 2 | 8 | 27.0 | 123.0 | 15.0 | |
| 10 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 2 | 33 | ... | 1 | 2 | 9 | 121.0 | 449.0 | 43.0 | |
| 11 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 2 | 34 | ... | 1 | 2 | 10 | 229.0 | 437.0 | 98.0 | 1 |

12 rows × 31 columns

# Part 1 – Get data into useable format

- Unzip data files
  - Each group of 30 files within its own folder
- Add in header data to each line
  - Add lat/long info based on day/time/location
- Read into database



| | julian_day | yearhour | month | day | hour | glob_horiz | dif_horiz | solar_altitude | solar_azimuth |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 8 | 1 | 1 | 8 | 7.0 | 6.0 | 5.44 | -53.16 |
| 1 | 1 | 9 | 1 | 1 | 9 | 10.0 | 21.0 | 13.83 | -42.02 |
| 2 | 1 | 10 | 1 | 1 | 10 | 69.0 | 65.0 | 20.46 | -29.34 |
| 3 | 1 | 11 | 1 | 1 | 11 | 173.0 | 173.0 | 24.77 | -15.12 |
| 4 | 1 | 12 | 1 | 1 | 12 | 122.0 | 121.0 | 26.26 | 0.09 |

```
In [24]: conn = sqlite3.connect("weather_database.db")
```

```
In [29]: df = pd.read_sql_query("SELECT * FROM weatherdata WHERE solar_altitude > 0 AND state = 'NJ' LIMIT 12;", conn)
```

```
In [30]: df
```

Out[30]:

| | station_id | city | state | timezone | lat | long | sm | elev | julian_day | yearhour | ... | month | day | hour | glob_horiz | dir_norm | dif_horiz | dir_horiz | skycover | solar_altitude | solar_azimuth |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 8 | ... | 1 | 1 | 8 | 7.0 | 5.0 | 6.0 | 1.0 | 10.0 | 5.37 | -53.30 |
| 1 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 9 | ... | 1 | 1 | 9 | 31.0 | 4.0 | 31.0 | 0.0 | 10.0 | 13.79 | -42.20 |
| 2 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 10 | ... | 1 | 1 | 10 | 68.0 | 6.0 | 66.0 | 2.0 | 10.0 | 20.46 | -29.54 |
| 3 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 11 | ... | 1 | 1 | 11 | 68.0 | 2.0 | 68.0 | 0.0 | 10.0 | 24.81 | -15.33 |
| 4 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 12 | ... | 1 | 1 | 12 | 89.0 | 7.0 | 86.0 | 3.0 | 10.0 | 26.34 | -0.11 |
| 5 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 13 | ... | 1 | 1 | 13 | 120.0 | 7.0 | 117.0 | 3.0 | 10.0 | 24.85 | 15.12 |
| 6 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 14 | ... | 1 | 1 | 14 | 83.0 | 5.0 | 81.0 | 2.0 | 10.0 | 20.54 | 29.34 |
| 7 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 15 | ... | 1 | 1 | 15 | 107.0 | 1.0 | 106.0 | 1.0 | 10.0 | 13.90 | 42.02 |
| 8 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 1 | 16 | ... | 1 | 1 | 16 | 53.0 | 1.0 | 53.0 | 0.0 | 10.0 | 5.50 | 53.15 |
| 9 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 2 | 32 | ... | 1 | 2 | 8 | 27.0 | 123.0 | 15.0 | 12.0 | 1.0 | 5.36 | -53.43 |
| 10 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 2 | 33 | ... | 1 | 2 | 9 | 121.0 | 449.0 | 43.0 | 78.0 | 1.0 | 13.80 | -42.33 |
| 11 | 14734 | NEWARK | NJ | -5.0 | 40.7 | 74.17 | 75.0 | 9.0 | 2 | 34 | ... | 1 | 2 | 10 | 229.0 | 437.0 | 98.0 | 131.0 | 4.0 | 20.49 | -29.67 |

12 rows × 21 columns

# CLASSICAL MACHINE LEARNING

Data is pre-categorized or numerical

Data is not labeled in any way

## SUPERVISED

## UNSUPERVISED

Predict a category

Predict a number

Divide by similarity

Identify sequences

### CLASSIFICATION
«Divide the socks by color»

### REGRESSION
«Divide the ties by length»

### CLUSTERING
«Split up similar clothing into stacks»

Find hidden dependencies

### ASSOCIATION
«Find what clothes I often wear together»

### DIMENSION REDUCTION (generalization)
«Make the best outfits from the given clothes»

From https://vas3k.com/blog/machine_learning/?ref=hn

# CLASSICAL MACHINE LEARNING

Data is pre-categorized
or numerical

Data is not labeled
in any way

## SUPERVISED

## UNSUPERVISED

Predict
a category

Predict
a number

Divide
by similarity

Identify sequences

## CLASSIFICATION

«Divide the socks by color»

## CLUSTERING

«Split up similar clothing
into stacks»

Find hidden
dependencies

## ASSOCIATION

«Find what clothes I often
wear together»

## REGRESSION

«Divide the ties by length»

## DIMENSION
REDUCTION
(generalization)

«Make the best outfits from the given clothes»

# Part 2 – Create machine learning model

# Part 2 – Create machine learning model



**A) Training**

Processed Data
**1961-1989**

Training Data
**80%**

Learning Algorithm

Model

**B) Validation**

Test Data
**20%**

Model

Validation

**C) Prediction**

New Data
**1990**

Model

Prediction

# Part 2 – Select train, validate model

# Results – Annual Hourly



Actual vs Predicted Diffuse Horizontal Illuminance

# Results – Annual Hourly – 12p.m. only



Actual vs Predicted Diffuse Horizontal Illuminance

# Analysis with Predicted Data - sDA

# Analysis with Predicted Data - Cumulative



Cumulative - Actual

Cumulative - Predicted

# But wait!
Aren't there already ways to do this?

# Existing Models

- Erbs et al., 1982 (ER)
- Orgill and Hollands, 1977 (OH)
- Reindl et al., 1990 (RE)
- Lam and Li, 1996 (LL)
- Skarteveit and Olseth, 1987 (SO)
- Louche et al., 1991 (LO)
- Maxwell, 1987 (MA)
- Vignola and McDaniels, 1984 (VM)

Sokol Dervishi and Ardeshir Mahdavi. Computing diffuse fraction of global horizontal solar radiation: A model comparison. Solar Energy, 2012

# Error Metrics

- Mean Bias Deviation (MBD)
- Relative Error (RE)
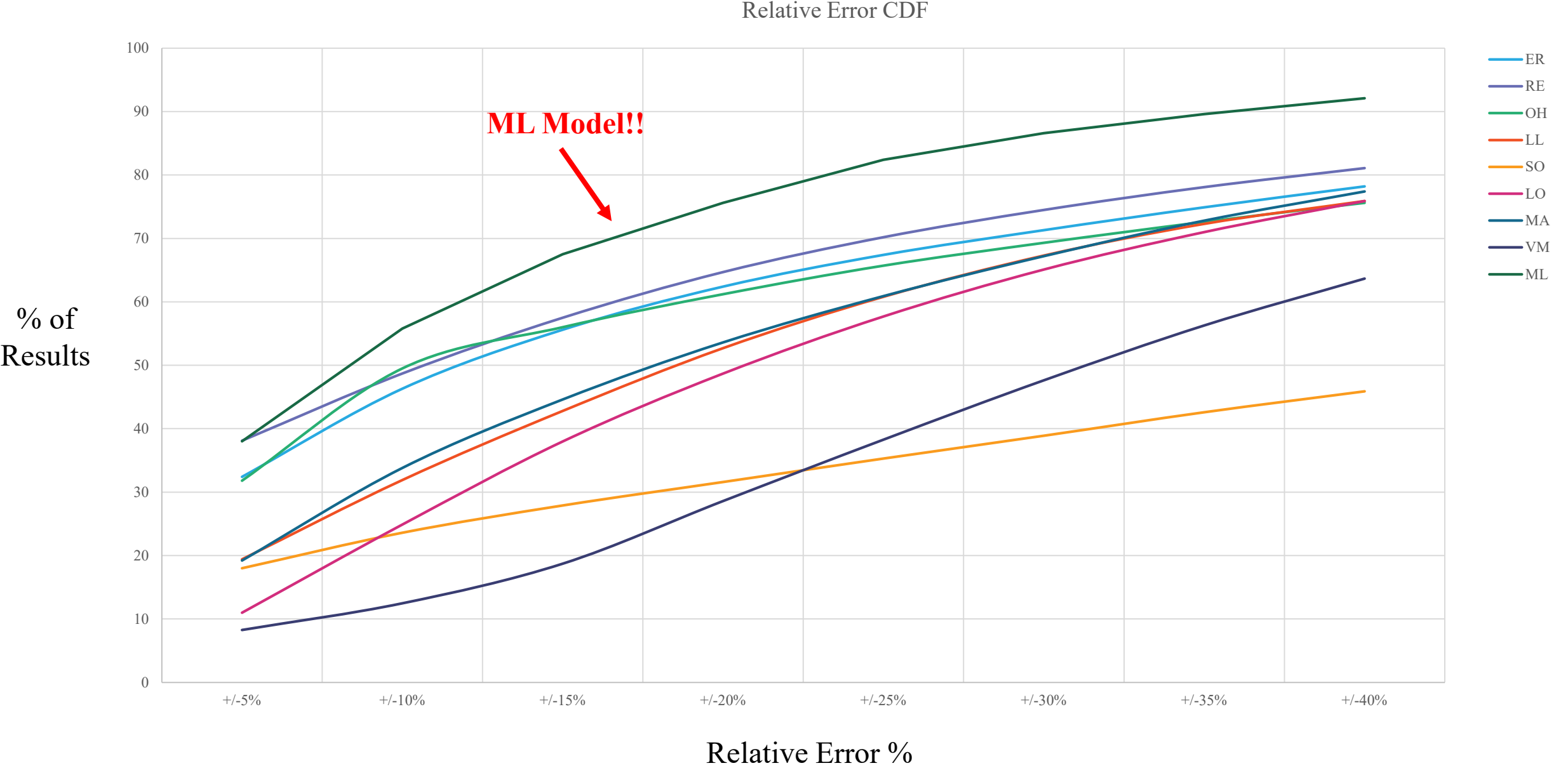- Root Mean Squared Deviation (RMSD, RMSE)

# Error Metrics – MBD and RMSD

| Model | MBD (%) | RMSD (W/m$^{-2}$) |
|---|---|---|
| ER | -9.2 | 37.4 |
| RE | -10.5 | 41.6 |
| OH | -13.3 | 43.1 |
| LL | 11.9 | 45.7 |
| SO | -98.3 | 199.9 |
| LO | 19.5 | 29.6 |
| MA | 21.1 | 33.2 |
| VM | -60.38 | 50.4 |

# Error Metrics – MBD and RMSD

| Model | MBD (%) | RMSD (W/m$^{-2}$) |
|-------|---------|-------------------|
| ER | -9.2 | 37.4 |
| RE | -10.5 | 41.6 |
| OH | -13.3 | 43.1 |
| LL | 11.9 | 45.7 |
| SO | -98.3 | 199.9 |
| LO | 19.5 | 29.6 |
| MA | 21.1 | 33.2 |
| VM | -60.38 | 50.4 |
| **ML** | **-5.87** | **35.32** |

# Error Metrics – Relative Error CDF



Relative Error CDF

# Takeaways

- Solution looks promising – needs more development
  - ML trained on only one NY weather station, ~131,000 measurements
  - 62M measurements in data set
  - Does using more than one weather station improve results, ie within a radius of target location where similar climate conditions are expected?
  - Does using all 236 weather stations improve results?
- Databases are very useful!
- Python/Jupyter environment worked well for this type of development.

# Thank you!

77 Water Street
New York, NY  10005

matt.franks@arup.com

+1 212 896 3000

ARUP