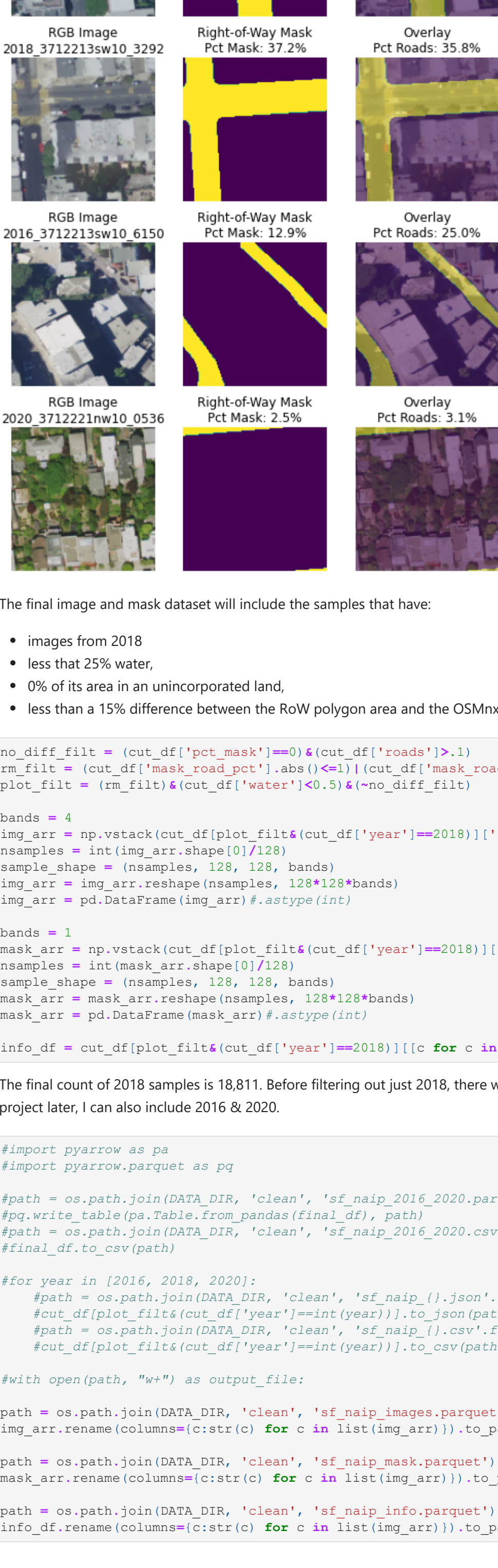



```
[41]: import matplotlib.pyplot as plt
cols, rows = 3, 5
axes = []
[axes.append((c,r)) for c in range(rows) for c in range(cols)]
axis_count = len(axes)
fig, ax = plt.subplots(rows, cols, figsize=(8.5, 16))

no_diff_filt = (cut_df['pot_mask']==0)&(cut_df['roads']>.1)
rm_filt = (cut_df['mask_road_pct'].abs()<=.1)|(cut_df['mask_road_diff'].abs()<=.15)
plot_filt = (rm_filt)&(cut_df['water']<0.5)&(no_diff_filt) &(cut_df['pot_roads']>.0)
plot_cols = ['id', 'image', 'mask', 'pot_mask', 'roads', 'mask_road_pct']

for row, (id, focus_image, focus_mask, pct_mask, pct_roads, pct_diff) in enumerate(cut_df[plot_filt].sample(n=100)):
    ax[row, 0].imshow(focus_image[:, :, :3])
    ax[row, 0].set_title('RGB Image' + '\n' + str(id)) #Size: (1x1x1).format(focus_image.shape[0], focus_image.shape[1], focus_image.shape[2])
    ax[row, 1].imshow(focus_mask)
    ax[row, 1].set_title('Right-of-Way Mask' + '\n' + 'Pot Mask: {0:.1%}'.format(pct_mask))
    ax[row, 2].imshow(focus_image[:, :, :3])
    ax[row, 2].imshow(focus_mask, alpha=.5)
    ax[row, 2].set_title('Pot Roads: {0:.1%}'.format(pct_roads) + '\n' + 'Pot Diff: {0:.1%}'.format(pct_diff))
    ax[row, 2].set_title('Overlay' + '\n' + 'Pot Roads: {0:.1%}'.format(pct_roads))

for col in range(cols):
    ax[row, col].set_axis_off()
```



The final image and mask dataset will include the samples that have:

- images from 2018
- less than 25% water.
- 0% of its area in an unincorporated land,
- less than a 15% difference between the RoW polygon area and the OSMnx centerline area

```
In [50]: no_diff_filt = (cut_df['pot_mask']==0)&(cut_df['roads']>.1)
rm_filt = (cut_df['mask_road_pct'].abs()<=.1)|(cut_df['mask_road_diff'].abs()<=.15)
plot_filt = (rm_filt)&(cut_df['water']<0.5)&(no_diff_filt) &(cut_df['pot_roads']>.0)

bands = 4
img_arr = np.vstack(cut_df[plot_filt&(cut_df['year']==2018)][['image']].values)
nsamples = int(img_arr.shape[0]/128)
sample_shape = (nsamples, 128, 128, bands)
img_arr = img_arr.reshape(nsamples, 128*128*bands)
img_arr = pd.DataFrame(img_arr).astype(int)

bands = 1
mask_arr = np.vstack(cut_df[plot_filt&(cut_df['year']==2018)][['mask']].values)
nsamples = int(mask_arr.shape[0]/128)
sample_shape = (nsamples, 128, 128, bands)
mask_arr = mask_arr.reshape(nsamples, 128*128*bands)
mask_arr = pd.DataFrame(mask_arr).astype(int)

info_df = cut_df[plot_filt&(cut_df['year']==2018)][[c for c in list(cut_df) if c not in ['image', 'mask']]].copy
```

The final count of 2018 samples is 18,811. Before filtering out just 2018, there were around 56,000 samples. So if I want to return to this project later, I can also include 2016 & 2020.

```
In [52]: #import pyarrow as pa
# import pyarrow.parquet as pq

#path = os.path.join(DATA_DIR, 'clean', 'sf_naisp_2016_2020.parquet')
#pq.write_table(pa.Table.from_pandas(final_df), path)
#path = os.path.join(DATA_DIR, 'clean', 'sf_naisp_2016_2020.csv')
#final_df.to_csv(path)

#for year in [2016, 2018, 2020]:
#    path = os.path.join(DATA_DIR, 'clean', 'sf_naisp_{}.json'.format(year))
#    cut_df[plot_filt&(cut_df['year']==int(year))].to_json(path)
#    path = os.path.join(DATA_DIR, 'clean', 'sf_naisp_{}.csv'.format(year))
#    cut_df[plot_filt&(cut_df['year']==int(year))].to_csv(path)

#with open(path, "w") as output_file:

path = os.path.join(DATA_DIR, 'clean', 'sf_naisp_images.parquet')
img_arr.rename(columns=[str(c) for c in list(img_arr)], to_parquet(path)

path = os.path.join(DATA_DIR, 'clean', 'sf_naisp_mask.parquet')
mask_arr.rename(columns=[str(c) for c in list(img_arr)], to_parquet(path)

path = os.path.join(DATA_DIR, 'clean', 'sf_naisp_info.parquet')
info_df.rename(columns=[str(c) for c in list(img_arr)], to_parquet(path)
```