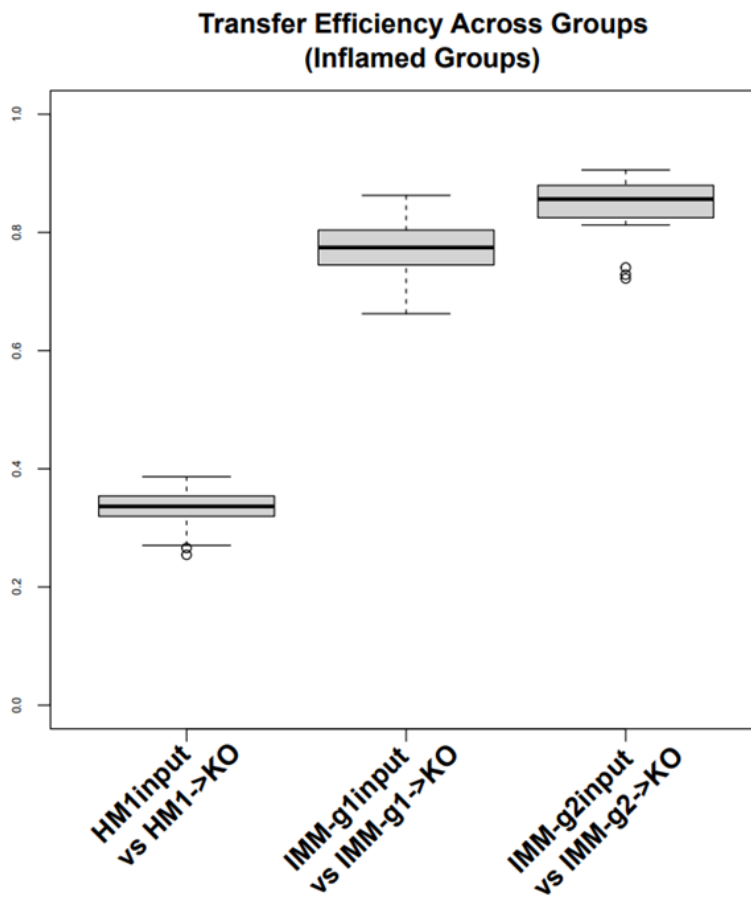


Contents

- Aim
- Main Python Methods
- Data Processing Steps
 - Step 1: Retrieve original counts table
 - Step 2: Retrieve original metadata table
 - Step 3: Create dictionary from metadata table
 - Step 4: Calculate total read counts per sample
 - Step 5: Extract counts for all six sample groups
 - Step 6: Normalize count values for all six sample groups
 - Step 7: Calculate Spearman Correlation Coefficients for each sample group
 - Step 8: Generate final figure
 - Step 9: Compare original and regenerated figures

Aim: Replicate Figure 4D



Main Python Methods

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import os
```

```

import re
from scipy import stats
from scipy.special import comb

# ***** MAIN METHODS *****
# Calculate Pearson correlation coefficients from input set of samples
def get_inter_corr_values(df1, df2):
    corr_values = [stats.pearsonr(df1.loc[:, col_A], df2.loc[:, col_B]).statistic \
                    for col_A in df1.columns for col_B in df2.columns]
    # Remove nan values
    corr_values = [val for val in corr_values if not np.isnan(val)]

    return corr_values

# Retrieve tab delimited file
def read_csv_file(file_path, skiprows=None, header = 0, sep = '\t', index_col=False):
    df = pd.read_csv(file_path, skiprows=skiprows, sep=sep, header=header, index_col=index_col)

    return df

# Create dictionary from metadata table
def get_dict_from_metadata(input_df):
    mydict = {}
    for row in input_df.iterrows():
        obj = row[1]
        sample_id = obj['SampleID']
        key = obj['FMTGroupFMTsourcegtRecipientbackground']
        if key not in mydict:
            mydict[key] = [sample_id]
        else:
            mydict[key].append(sample_id)
    return mydict

# ***** HELPER METHODS *****
**
# Normalize values in counts table
def _get_norm_counts(input_df, ser_sample_count_sums):
    overall_mean_count = ser_sample_count_sums.mean()
    df_norm_logged = pd.DataFrame()
    for col_name in input_df.columns:
        df_norm_logged.loc[:, col_name] = \
            np.log10(((input_df.loc[:, col_name] / ser_sample_count_sums[col_name]
            ]) * overall_mean_count) + 1)

    return df_norm_logged

# ***** INPUT FILE PATHS *****
****
current_working_dir = os.getcwd()
base_dir = r'C:\Users\young\Documents\anh_validation'
# original counts table
asv_tbl_file_path = os.path.join(base_dir, 'asv_biom-with-taxonomy.txt')
# original metadata table
metadata_file_path = os.path.join(base_dir, 'mappingMetadata.txt')

```

Data Processing Steps

Step 1: Retrieve original counts table

```
In [2]: df_asv = read_csv_file(asv_tbl_file_path, 1)
df_asv = df_asv.astype({col:'int32' for col in df_asv.columns[1:-1] }, copy=False)
df_asv
```

Out[2]:

	#OTU ID	1gKO.1	1gKO.2	1gKO.3	1gWT.1	1gWT.2	1gWT.3	2gKO.1
0	1ba8c796d07406783c96d016a6a5cace	13615	16637	17148	20227	23630	25656	14832
1	a6c38249aff7768283faf6cfbdeb05a8	26439	30129	19743	8955	10759	7074	18489
2	062f38ff92cfaee0654200b6f5be5ddf	7451	8774	8754	174	214	148	21958
3	1183cc23f552d81e63c93ca9fcb2f2c	225	223	184	13762	16856	18692	269
4	5e15ecfb579e72bf87c0bea3920bbf42	10108	12117	8633	10027	11910	7424	5979
...
4070	92bb8f4683ef5c8651e7d34dbb37ab2e	0	0	0	0	0	0	0
4071	92f09070a4fd5786bb34e756217e6ee1	0	0	0	0	0	0	0
4072	919b82324c41ed0046323c63aa1550da	0	0	0	0	0	0	0
4073	dbc0dad15ec1c8ad9d826cab94e18696	0	0	0	0	0	0	0
4074	1ff2d07d10264c23dc43e08d3097cd7c	0	0	0	0	0	0	0

4075 rows × 112 columns

Step 2: Retrieve original metadata table

```
In [3]: df_metadata = read_csv_file(metadata_file_path)
df_metadata
```

Out[3]:

	SampleID	UniversalCageNumber	Background	FMTGroupFMTsourcecgtRecipientbackground	Passage
0	F8-1	F8-cage-1	129.IL10KO	1gKOgtKO	8
1	F8-2	F8-cage-1	129.IL10KO	1gKOgtKO	8
2	F8-3	F8-cage-2	129.IL10KO	1gKOgtKO	8
3	F8-4	F8-cage-2	129.IL10KO	1gKOgtKO	8
4	F8-5	F8-cage-3	129.IL10KO	1gKOgtKO	8
...

105	1gWT.2	NaN	NaN	1gWTinput	1gWT
106	1gWT.3	NaN	NaN	1gWTinput	1gWT
107	2gWT.1	NaN	NaN	2gWTinput	2gWT
108	2gWT.2	NaN	NaN	2gWTinput	2gWT
109	2gWT.3	NaN	NaN	2gWTinput	2gWT

110 rows × 8 columns

Step 3: Create dictionary from metadata table

```
In [4]: dict_metadata = get_dict_from_metadata(df_metadata)
dict_metadata.keys()
```

```
Out[4]: dict_keys(['1gKOgtKO', '2gKOgtKO', '1gWTgtKO', '1gWTgtWT', '2gWTgtWT', 'hFMT.1.2.3.gtKO', 'hFMT.3.4.5.gtKO', 'hFMT.1.2.3.gtWT', 'hFMT.1.2.3.input', 'hFMT.3.4.5.input', '1gKOinput', '2gKOinput', '1gWTinput', '2gWTinput'])
```

Step 4: Calculate total read counts per sample

```
In [5]: sample_count_sums = df_asv.iloc[:, 1:-1].sum(axis=0)
sample_count_sums
```

```
Out[5]: 1gKO.1      118256
1gKO.2      141891
1gKO.3      123292
1gWT.1      119717
1gWT.2      146158
...
h1-2-3.2    135326
h1-2-3.3    133745
h3-4-5.1    129613
h3-4-5.2    140316
h3-4-5.3    132984
Length: 110, dtype: int64
```

Step 5: Extract counts for target sample groups

```
In [6]: # 'hFMT.1.2.3.input' --> 'HM1Input'
key_name = 'hFMT.1.2.3.input'
print('Grouped columns: {}'.format(dict_metadata[key_name]))
HM1_Input = df_asv[[col_name for col_name in dict_metadata[key_name] ]]
HM1_Input
```

Grouped columns: ['h1-2-3.1', 'h1-2-3.2', 'h1-2-3.3']

Out[6]:

	h1-2-3.1	h1-2-3.2	h1-2-3.3
0	1483	1515	1402
1	1129	691	725
2	2927	3523	3254
3	298	231	238

4	218	160	111
...
4070	0	0	0
4071	0	0	0
4072	0	0	0
4073	0	0	0
4074	0	0	0

4075 rows × 3 columns

```
In [7]: # 'hFMT.1.2.3.gtKO' --> 'HM1->KO'
key_name = 'hFMT.1.2.3.gtKO'
print('Grouped columns: {}'.format(dict_metadata[key_name]))
HM1_KO = df_asv[[col_name for col_name in dict_metadata[key_name] ]]
HM1_KO
```

Grouped columns: ['F3-7', 'F3-8', 'F3-9', 'F3-10', 'F3-11', 'F3-12', 'F1-7', 'F1-8', 'F1-9', 'F1-10', 'F1-11', 'F1-12', 'F1-13', 'F1-14', 'F1-15']

Out[7]:

	F3-7	F3-8	F3-9	F3-10	F3-11	F3-12	F1-7	F1-8	F1-9	F1-10	F1-11	F1-12	F1-13	F
0	10097	574	12959	18737	17857	12650	10738	9576	7304	14409	14382	8787	7694	
1	35102	51633	31840	22887	20446	33188	23565	21707	12802	26719	5701	15715	25125	
2	1328	91	2203	189	182	244	33299	13739	53077	17570	34094	18922	19204	4
3	137	315	305	134	221	290	203	348	363	238	219	205	255	
4	34960	13573	4109	3811	1155	28938	1144	6973	666	8230	1287	11346	13195	
...	
4070	0	0	0	0	0	0	0	0	0	0	0	0	0	
4071	0	0	0	0	0	0	0	0	0	0	0	0	0	
4072	0	0	0	0	0	0	0	0	0	0	0	0	0	
4073	0	0	0	0	0	0	0	0	0	0	0	0	0	
4074	0	0	0	0	0	0	0	0	0	0	0	0	0	

4075 rows × 15 columns

```
In [8]: # '1gKOinput' --> 'IMM-glinput'
key_name = '1gKOinput'
print('Grouped columns: {}'.format(dict_metadata[key_name]))
IMM_glinput = df_asv[[col_name for col_name in dict_metadata[key_name] ]]
IMM_glinput
```

Grouped columns: ['1gKO.1', '1gKO.2', '1gKO.3']

Out[8]:

	1gKO.1	1gKO.2	1gKO.3
0	13615	16637	17148
1	26439	30129	19743
2	7451	8774	8754
3	225	223	184
4	10108	12117	8633

4	10100	12117	0033
...
4070	0	0	0
4071	0	0	0
4072	0	0	0
4073	0	0	0
4074	0	0	0

4075 rows × 3 columns

```
In [9]: # '1gKOgtKO' --> 'IMM-g1->KO'
key_name = '1gKOgtKO'
print('Grouped columns: {}'.format(dict_metadata[key_name]))
IMM_g1_KO = df_asv[[col_name for col_name in dict_metadata[key_name] ]]
IMM_g1_KO
```

Grouped columns: ['F8-1', 'F8-2', 'F8-3', 'F8-4', 'F8-5', 'F8-6', 'F8-7', 'F8-8', 'F4-1', 'F4-2', 'F4-3', 'F4-4', 'F4-5', 'F3-1', 'F3-2', 'F3-4', 'F3-5', 'F3-6']

Out[9]:

	F8-1	F8-2	F8-3	F8-4	F8-5	F8-6	F8-7	F8-8	F4-1	F4-2	F4-3	F4-4	F4-5	
0	3339	15223	11423	18205	12486	10761	9339	466	29629	20134	27588	26741	25644	2
1	10283	6602	9314	13651	8789	16118	12045	13841	3687	6998	6904	7796	4011	
2	41532	40602	52337	11038	16781	8926	25758	31166	1627	37042	23825	26508	30376	
3	240	276	331	302	341	231	226	306	275	198	213	271	172	2
4	1806	6778	6725	12734	3901	10948	10347	1234	8891	3718	9188	4419	8315	
...	
4070	0	0	0	0	0	0	0	0	0	0	0	0	0	
4071	0	0	0	0	0	0	0	0	0	0	0	0	0	
4072	0	0	0	0	0	0	0	0	0	0	0	0	0	
4073	0	0	0	0	0	0	0	0	0	0	0	0	0	
4074	0	0	0	0	0	0	0	0	0	0	0	0	0	

4075 rows × 18 columns

```
In [10]: # '2gKOinput' --> 'IMM-g2input'
key_name = '2gKOinput'
print('Grouped columns: {}'.format(dict_metadata[key_name]))
IMM_g2input = df_asv[[col_name for col_name in dict_metadata[key_name] ]]
IMM_g2input
```

Grouped columns: ['2gKO.1', '2gKO.2', '2gKO.3']

Out[10]:

	2gKO.1	2gKO.2	2gKO.3
0	14832	14973	13992
1	18489	18178	9122
2	21958	22435	17629
3	269	153	230
4	5979	4360	2762

4	3373	4330	2702
...
4070	0	0	0
4071	0	0	0
4072	0	0	0
4073	0	0	0
4074	0	0	0

4075 rows × 3 columns

```
In [11]: # '2gKOgtKO' --> 'IMM-g2->KO'
key_name = '2gKOgtKO'
print('Grouped columns: {}'.format(dict_metadata[key_name]))
IMM_g2_KO = df_asv[[col_name for col_name in dict_metadata[key_name] ]]
IMM_g2_KO
```

Grouped columns: ['F8-9', 'F8-10', 'F8-11', 'F8-12', 'F8-13', 'F8-14', 'F8-15', 'F8-16']

Out[11]:

	F8-9	F8-10	F8-11	F8-12	F8-13	F8-14	F8-15	F8-16
0	23451	24742	31312	13542	4593	10228	14477	23819
1	8030	14203	27872	26539	12663	13720	23316	14430
2	19653	22616	4603	28535	20096	20133	14650	24062
3	266	344	300	269	257	330	249	226
4	6490	6143	12353	2229	6367	11645	9607	10622
...
4070	0	0	0	0	0	0	0	0
4071	0	0	0	0	0	0	0	0
4072	0	0	0	0	0	0	0	0
4073	0	0	0	0	0	0	0	0
4074	0	0	0	0	0	0	0	0

4075 rows × 8 columns

Step 6: Normalize count values for target sample groups

```
In [12]: HM1_Input_norm = _get_norm_counts(HM1_Input, sample_count_sums)
HM1_Input_norm
```

Out[12]:

	h1-2-3.1	h1-2-3.2	h1-2-3.3
0	3.170992	3.168281	3.139741
1	3.052636	2.827698	2.853624
2	3.466129	3.534613	3.505231
3	2.475230	2.353116	2.371106
4	2.340004	2.194478	2.041962

...
4070	0.000000	0.000000	0.000000
4071	0.000000	0.000000	0.000000
4072	0.000000	0.000000	0.000000
4073	0.000000	0.000000	0.000000
4074	0.000000	0.000000	0.000000

4075 rows × 3 columns

```
In [13]: HM1_KO_norm = _get_norm_counts(HM1_KO, sample_count_sums)
HM1_KO_norm
```

Out[13]:

	F3-7	F3-8	F3-9	F3-10	F3-11	F3-12	F1-7	F1-8	F1-9	F1-
0	4.023995	2.715632	4.028643	4.300323	4.281648	4.108629	4.026117	3.977907	3.888433	4.1297
1	4.565105	4.668820	4.419021	4.387208	4.340445	4.527498	4.367439	4.333297	4.132126	4.3978
2	3.143272	1.920176	3.259285	2.306214	2.292111	2.395644	4.517597	4.134664	4.749730	4.2158
3	2.159501	2.455717	2.402055	2.157737	2.376040	2.470375	2.304807	2.539512	2.585851	2.3495
4	4.563345	4.088595	3.529896	3.608747	3.092750	4.467988	3.053962	3.840159	2.848906	3.8865
...
4070	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0000
4071	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0000
4072	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0000
4073	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0000
4074	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0000

4075 rows × 15 columns

```
In [14]: IMM_glinput_norm = _get_norm_counts(IMM_glinput, sample_count_sums)
IMM_glinput_norm
```

Out[14]:

	1gKO.1	1gKO.2	1gKO.3
0	4.180178	4.188104	4.262257
1	4.468392	4.446001	4.323454
2	3.918399	3.910251	3.970273
3	2.400047	2.317401	2.295045
4	4.050836	4.050434	3.964229
...
4070	0.000000	0.000000	0.000000
4071	0.000000	0.000000	0.000000
4072	0.000000	0.000000	0.000000
4073	0.000000	0.000000	0.000000
4074	0.000000	0.000000	0.000000

4075 rows × 3 columns


```
In [15]: IMM_g1_KO_norm = _get_norm_counts(IMM_g1_KO, sample_count_sums)
IMM_g1_KO_norm
```

Out[15]:

	F8-1	F8-2	F8-3	F8-4	F8-5	F8-6	F8-7	F8-8	F4-1	F4
0	3.530045	4.171099	4.067413	4.239122	4.320109	4.065412	3.978399	2.753708	4.509011	4.2994
1	4.018461	3.808313	3.978778	4.114104	4.167634	4.240858	4.088895	4.225749	3.604062	3.8405
2	4.624693	4.597128	4.728413	4.021838	4.448499	3.984224	4.418981	4.578248	3.248913	4.5642
3	2.388291	2.431092	2.530706	2.460420	2.757180	2.398871	2.364043	2.571443	2.478060	2.2944
4	3.263255	3.819737	3.837351	4.083907	3.814908	4.072893	4.022908	3.176160	3.986276	3.5659
...
4070	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0000
4071	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0000
4072	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0000
4073	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0000
4074	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0000

4075 rows × 18 columns

```
In [16]: IMM_g2input_norm = _get_norm_counts(IMM_g2input, sample_count_sums)
IMM_g2input_norm
```

Out[16]:

	2gKO.1	2gKO.2	2gKO.3
0	4.143317	4.154522	4.242762
1	4.239025	4.238754	4.056985
2	4.313700	4.330129	4.343104
3	2.403557	2.166842	2.460093
4	3.748792	3.618774	3.538206
...
4070	0.000000	0.000000	0.000000
4071	0.000000	0.000000	0.000000
4072	0.000000	0.000000	0.000000
4073	0.000000	0.000000	0.000000
4074	0.000000	0.000000	0.000000

4075 rows × 3 columns

```
In [17]: IMM_g2_KO_norm = _get_norm_counts(IMM_g2_KO, sample_count_sums)
IMM_g2_KO_norm
```

Out[17]:

	F8-9	F8-10	F8-11	F8-12	F8-13	F8-14	F8-15	F8-16
0	4.480880	4.378895	4.379459	4.143025	3.698748	4.105404	4.195335	4.353703
1	4.015462	4.137854	4.328918	4.435212	4.139133	4.232959	4.402300	4.136058
2	4.404150	4.339877	3.546894	4.466704	4.339694	4.399505	4.200494	4.358111

3	2.536849	2.523304	2.362739	2.442634	2.448049	2.615148	2.432436	2.332891
4	3.923001	3.773895	3.975548	3.359611	3.840562	4.161749	4.017259	4.003009
...
4070	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
4071	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
4072	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
4073	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
4074	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000

4075 rows × 8 columns

Step 7: Calculate Spearman Correlation Coefficients for each sample group

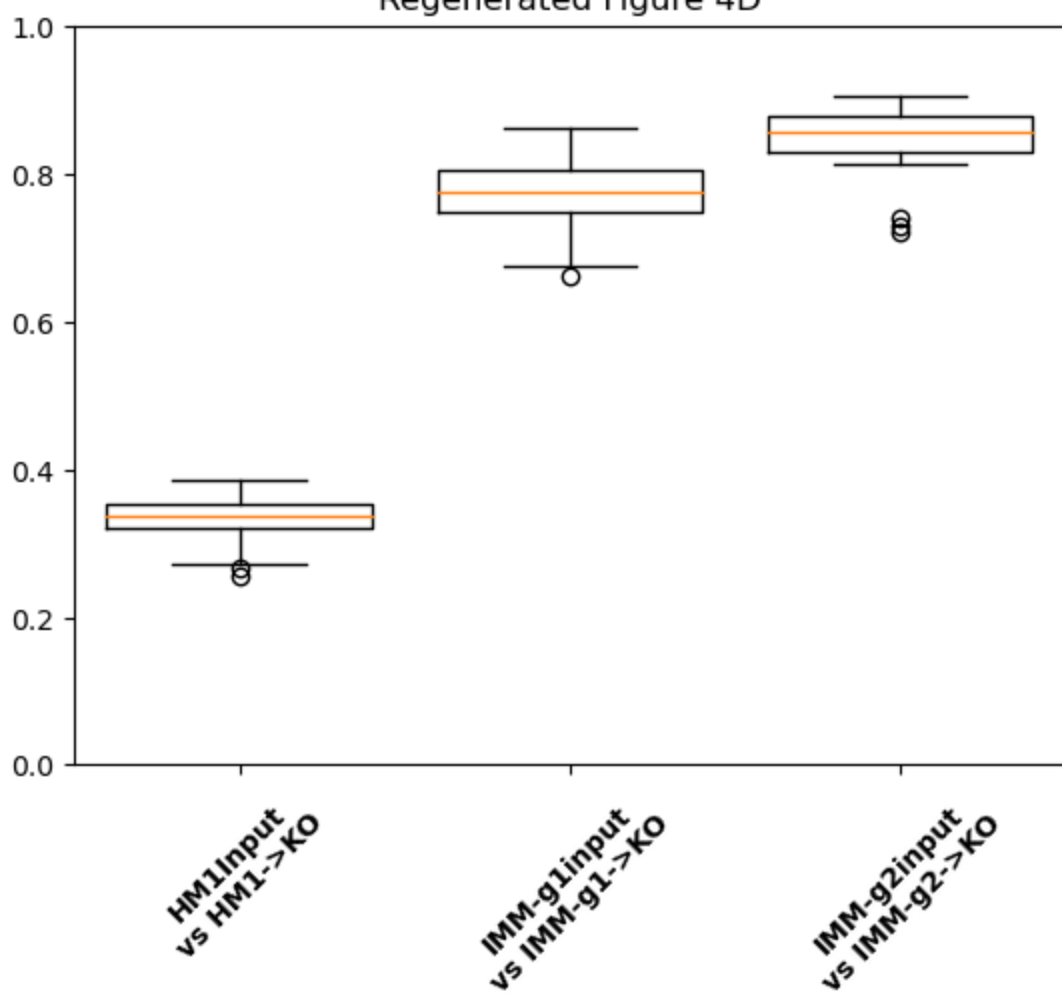
```
In [18]: corr_value1 = get_inter_corr_values(HM1_Input_norm, HM1_KO_norm)
corr_value2 = get_inter_corr_values(IMM_g1input_norm, IMM_g1_KO_norm)
corr_value3 = get_inter_corr_values(IMM_g2input_norm, IMM_g2_KO_norm)
```

Step 8: Generate final figure

```
In [19]: data = [corr_value1, corr_value2, corr_value3]
x_tick_labels = ['HM1Input\nvs HM1->KO',
                  'IMM-g1input\nvs IMM-g1->KO',
                  'IMM-g2input\nvs IMM-g2->KO']

fig, ax = plt.subplots()
ax.boxplot(data, widths=0.8)
ax.set_title('Regenerated Figure 4D')
ax.set_yticks([0.0, 0.2, 0.4, 0.6, 0.8, 1.0])
ax.set_xticklabels(x_tick_labels,
                  rotation=45, fontsize=10, weight='bold')

plt.show()
```



Step 9: Compare original and regenerated figures

