

```
| : # Get better looking pictures
| : %config InlineBackend.figure_format = 'retina'
|
| : import pandas as pd
```

```

import scipy.stats as st
import matplotlib.pyplot as plt
import datetime
from dateutil import parser
import networkx as nx
import random as r

```

```
In [3]:
df = pd.read_feather('enron.feather')
df = df.sort_values(['Date'])
df.tail(5)
```

	MailID	Date	From	To	Recipients	Subject	filename
601933	72015	2002-07-12	denise.williams	ge_benefits	1	URGENT!!! CUTOVER WEEKEND	fischer-m/all_documents/428.
601861	71978	2002-07-12	mark.fisher	tom.nemila	1	WR613 Pitch System performance	fischer-m/all_documents/425.
601579	71905	2002-07-12	mark.fisher	tom.nemila	1	WR613 Pitch System performance	fischer-m/_sent_mail/1.
605086	73118	2002-07-12	denise.williams	ge_benefits	1	URGENT!!! CUTOVER WEEKEND	fischer-m/notes_inbox/2.
601581	71907	2002-07-12	mark.fisher	tom.nemila	1	WR627 Fault Paretos (May 2002)	fischer-m/_sent_mail/2.

Email traffic over time

Group the data set by `Date` and `MailID`, which will get you an index that collects all of the unique mail IDs per date. Then reset the index so that those date and mail identifiers become columns and then select for just those columns; we don't actually care about the counts created by the `groupby` (that was just to get the index). Create a histogram that shows the amount of traffic per day. Then specifically for email sent from `richard.shapiro` and then `john.lavorato`. Because some dates are set improperly (to 1980), filter for dates greater than January 1, 1999.

```
In [4]:
stored_df=df.groupby(["Date","MailID"])['Recipients'].nunique().reset_index(drop=False)
del stored_df["Recipients"]
stored_df_richard=df[df["From"]=="richard.shapiro"].groupby(["Date","MailID"])['Recipients'].nunique().reset_index()
stored_df_johnl=df[df["From"]=="john.lavorato"].groupby(["Date","MailID"])['Recipients'].nunique().reset_index()
del stored_df_richard["Recipients"]
del stored_df_johnl["Recipients"]
```

```
In [5]:
fig, ax = plt.subplots(figsize=(12,4))
ax.hist(stored_df[stored_df["Date"]>"1999-01-01"]["Date"], bins=100, color="#fee090")
ax.set_xlabel("Date")
ax.set_ylabel("Number of emails")
ax.set_title("Histogram of emails sent")

ax.tick_params(axis='x', rotation=45)

for rect in ax.patches:
    rect.set_linewidth(.5)
    rect.set_edgecolor('grey')

ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.spines['left'].set_linewidth(0.5)
ax.spines['bottom'].set_linewidth(0.5)

plt.show()
```

```
In [6]:
fig, ax = plt.subplots(figsize=(12,4))
ax.hist(stored_df_richard[stored_df_richard["Date"]>"1999-01-01"]["Date"], bins=100, color="#fee090")
ax.set_xlabel("Date")
ax.set_ylabel("Number of emails")
ax.set_title("Histogram of emails sent by richard.shapiro")

ax.tick_params(axis='x', rotation=45)

for rect in ax.patches:
    rect.set_linewidth(.5)
    rect.set_edgecolor('grey')

ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.spines['left'].set_linewidth(0.5)
ax.spines['bottom'].set_linewidth(0.5)

plt.show()
```

Year	Cases
1999-10	5
2000-01	7
2001-02	7
2002-03	7
2003-04	18
2004-05	10
2005-06	11

```
fig, ax = plt.subplots(figsize=(12,4))
ax.hist(stored_df_johnl[stored_df_johnl["Date"]>"1999-01-01"]["Date"], bins=100, color="#fee090")
ax.set_xlabel("Date")
ax.set_ylabel("Number of emails")
ax.set_title("Hisotgram of emails sent by john.lavorato")

ax.tick_params(axis='x', rotation=45)

for rect in ax.patches:
    rect.set_linewidth(.5)
    rect.set_edgecolor('grey')

ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.spines['left'].set_linewidth(0.5)
ax.spines['bottom'].set_linewidth(0.5)

plt.show()
```

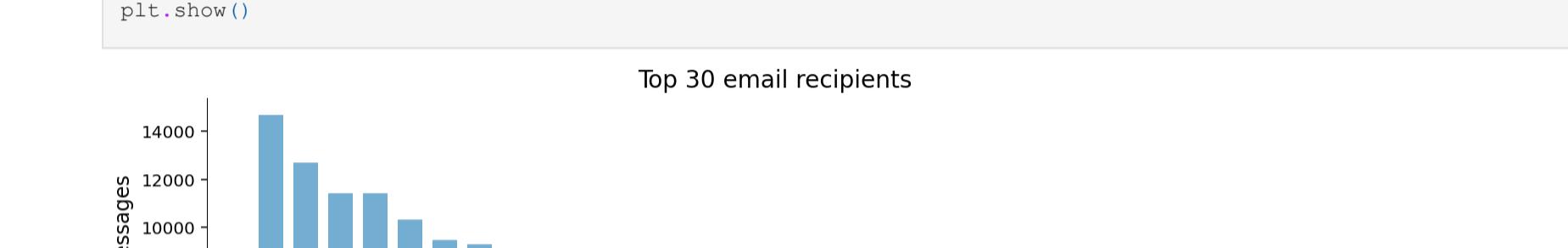
Hisotgram of emails sent by john.lavorato

A histogram titled "Hisotgram of emails sent by john.lavorato". The x-axis is labeled "Date" and shows a range from 1999-01-01 to the present. The y-axis is labeled "Number of emails" and ranges from 120 to 160. A single bar is visible at the start of the x-axis, representing the count of emails sent on January 1, 1999, which is approximately 165.

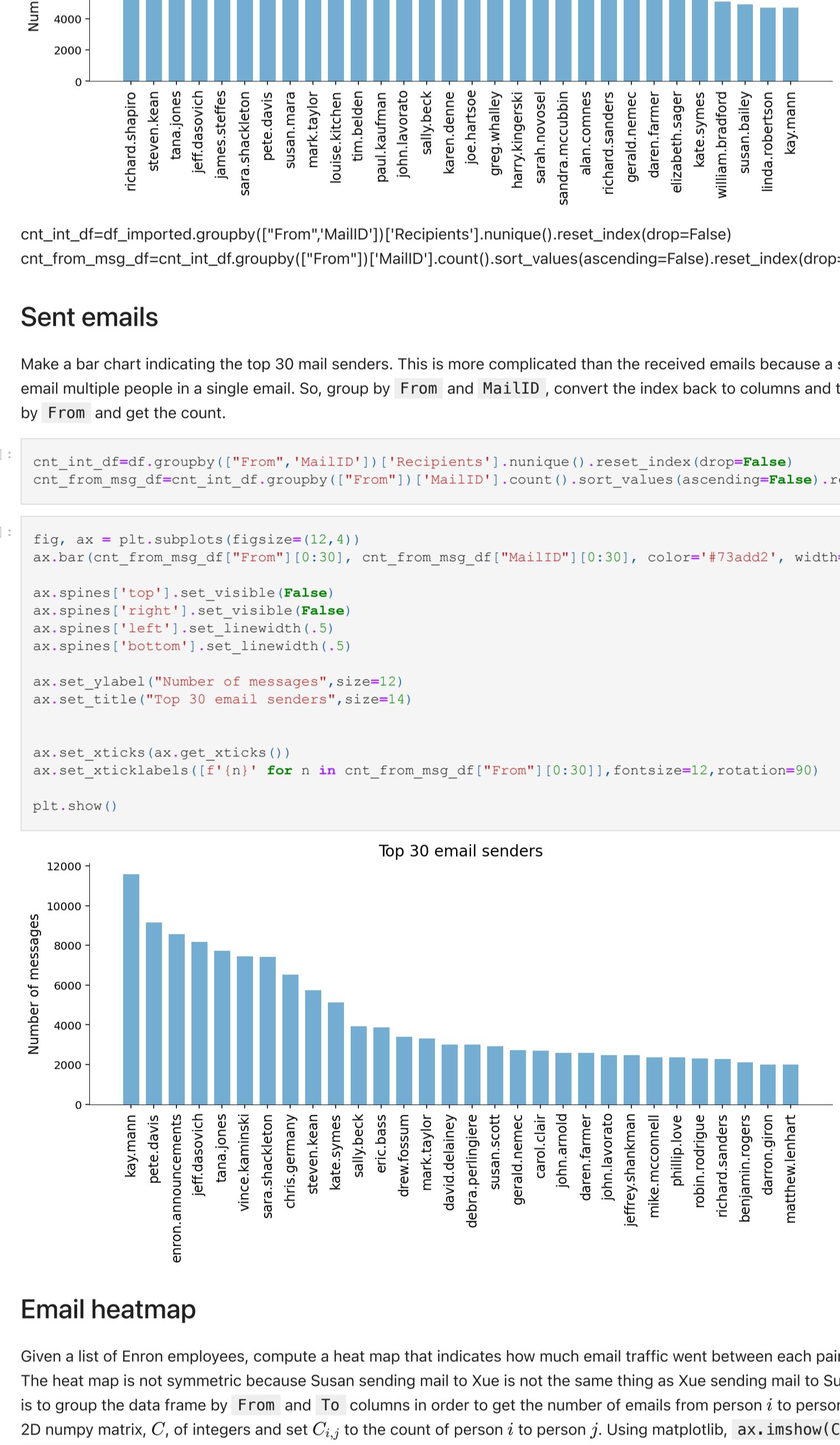
A bar chart titled "Received emails" showing the count of messages received per user over time. The y-axis is labeled "Number of messages" and ranges from 0 to 80. The x-axis is labeled "Date" and shows monthly intervals from January 2000 to February 2002. The bars are colored yellow and orange, representing different users. The chart shows a peak in activity around July 2000, followed by a decline and then a smaller peak around October 2001.

Date	User 1 (Yellow)	User 2 (Orange)
2000-01	35	10
2000-02	65	15
2000-03	25	5
2000-04	25	18
2000-05	35	22
2000-06	38	28
2000-07	80	65
2000-08	85	45
2000-09	75	55
2000-10	55	58
2000-11	78	50
2000-12	82	45
2001-01	45	48
2001-02	20	42
2001-03	5	3
2001-04	10	2
2001-05	22	1
2001-06	18	0
2001-07	5	0
2001-08	3	0
2001-09	10	0
2001-10	25	5
2001-11	30	10
2001-12	15	8
2002-01	12	15
2002-02	10	12

```
In [9]:  
fig, ax = plt.subplots(figsize=(12, 4))  
ax.bar(cnt_messages_df["To"][:30], cnt_messages_df["MailID"][:30], color="#73add2", width=0.7)  
  
ax.spines['top'].set_visible(False)  
ax.spines['right'].set_visible(False)  
ax.spines['left'].set_linewidth(.5)  
ax.spines['bottom'].set_linewidth(.5)  
  
ax.set_ylabel("Number of messages", size=12)  
ax.set_title("Top 30 email recipients", size=14)  
  
ax.set_xticks(ax.get_xticks())  
ax.set_xticklabels([f'{n}' for n in cnt_messages_df["To"][:30]], fontsize=12, rotation=90)  
  
plt.show()
```



Recipient Index	Number of messages
0	14500
1	12800
2	11500
3	11500
4	10500
5	9800
6	9500
7	9200
8	9000
9	8800
10	8500
11	8300
12	8100
13	7900
14	7700
15	7500
16	7300
17	7100
18	6900
19	6700
20	6500
21	6300
22	6100
23	5900
24	5700
25	5500
26	5300
27	5100
28	4900
29	4700



```
ax.text(), the coordinates are X,Y whereas the coordinates in the C matrix are row,column so you will have to flip the
coordinates.

In [12]:
people = ['jeff.skilling', 'kenneth.lay', 'louise.kitchen', 'tana.jones',
          'sara.shackleton', 'vince.kaminski', 'sally.beck', 'john.lavorato',
          'mark.taylor', 'greg.whalley', 'jeff.dasovich', 'steven.kean',
          'chris.germany', 'mike.mcconnell', 'benjamin.rogers', 'j.kaminski',
          'stanley.horton', 'a..shankman', 'richard.shapiro']

In [13]:
cnt_to_frm_df=df.groupby(["To",'From'])['MailID'].count().reset_index(drop=False)
cnt_to_frm_req_df=cnt_to_frm_df[cnt_to_frm_df['To'].isin(people) & cnt_to_frm_df['From'].isin(people)]
elem_df=cnt_to_frm_req_df.pivot(index='From', columns='To', values='MailID')
elem_df=elem_df[people]
elem_df=elem_df.reindex(people).fillna(0)

In [14]:
C=elem_df.to_numpy()
fig, ax = plt.subplots(figsize=(200, 8))

im = ax.imshow(C, cmap='GnBu', vmax=4000)

ax.set_xticks(range(0,19))
```

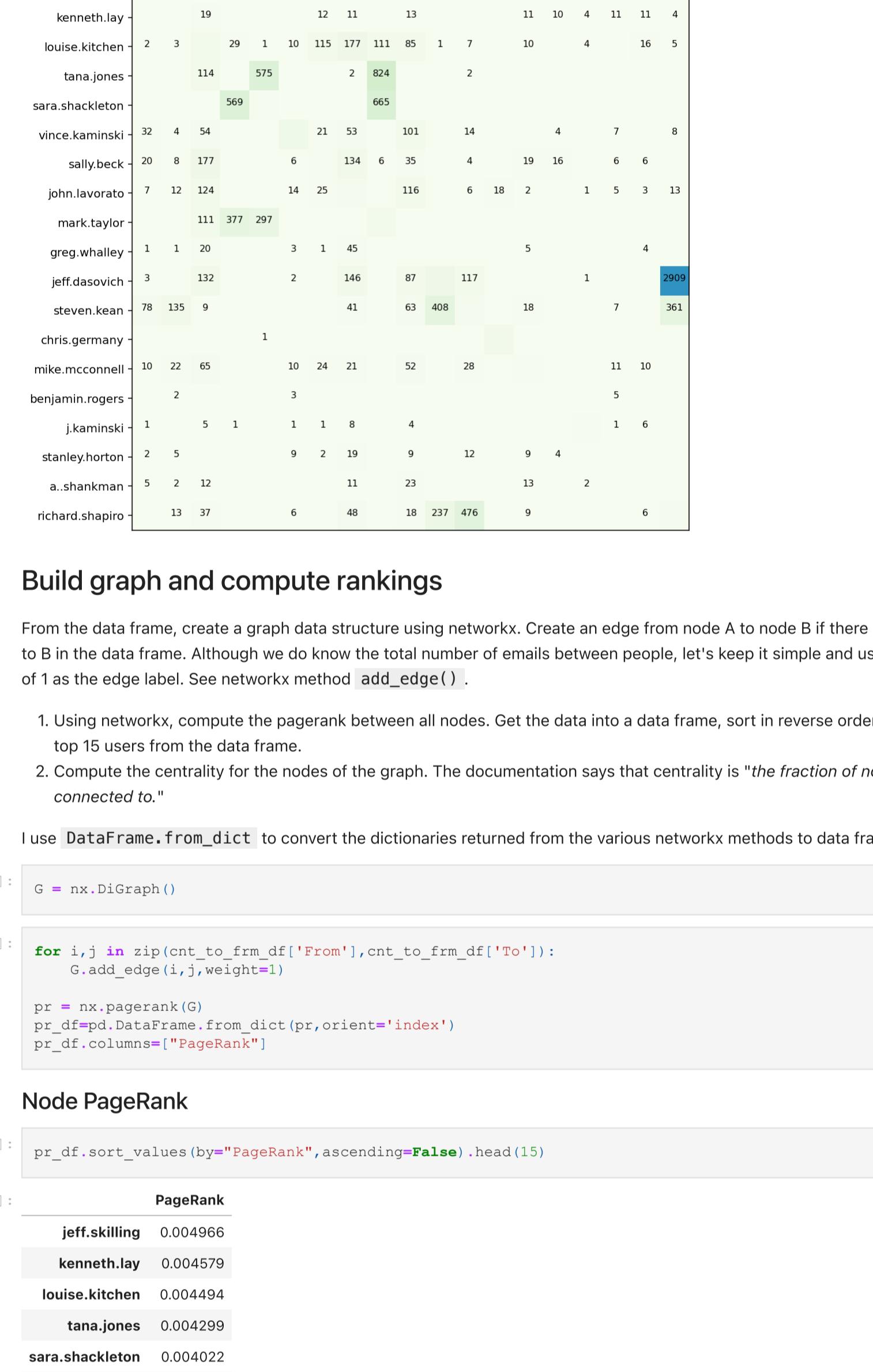
```
ax.set_yticks(range(0,19))
ax.set_yticklabels(people)

for i in range(19):
    for j in range(19):
        if i!=j:
            if C[j,i]!=0:
                ax.text(i, j, f"{C[j,i]:.0f}", horizontalalignment='center',size=8)

ax.xaxis.tick_top()

#fig.colorbar(im, ax=ax, shrink=.83)
plt.tight_layout()
plt.show()
```





john.lavorato	0.00
gerald.nemec	0.00
rod.hayslett	0.00
mark.taylor	0.00
greg.whalley	0.00

daren.farmer	0.002007
steven.kean	0.002003

Out[18]:

PageRank

sally.beck 0.088048

outlook.team 0.082401

david.forster 0.079988

kenneth.lay 0.076137

technology.enron 0.063662

jeff.skilling 0.058220

tana.jones 0.054677

louise.kitchen 0.053445

jeff.dasovich 0.048773

sara.shackleton 0.047952

julie.clyatt	0.045847
bodyshop	0.044922
david.oxley	0.043793

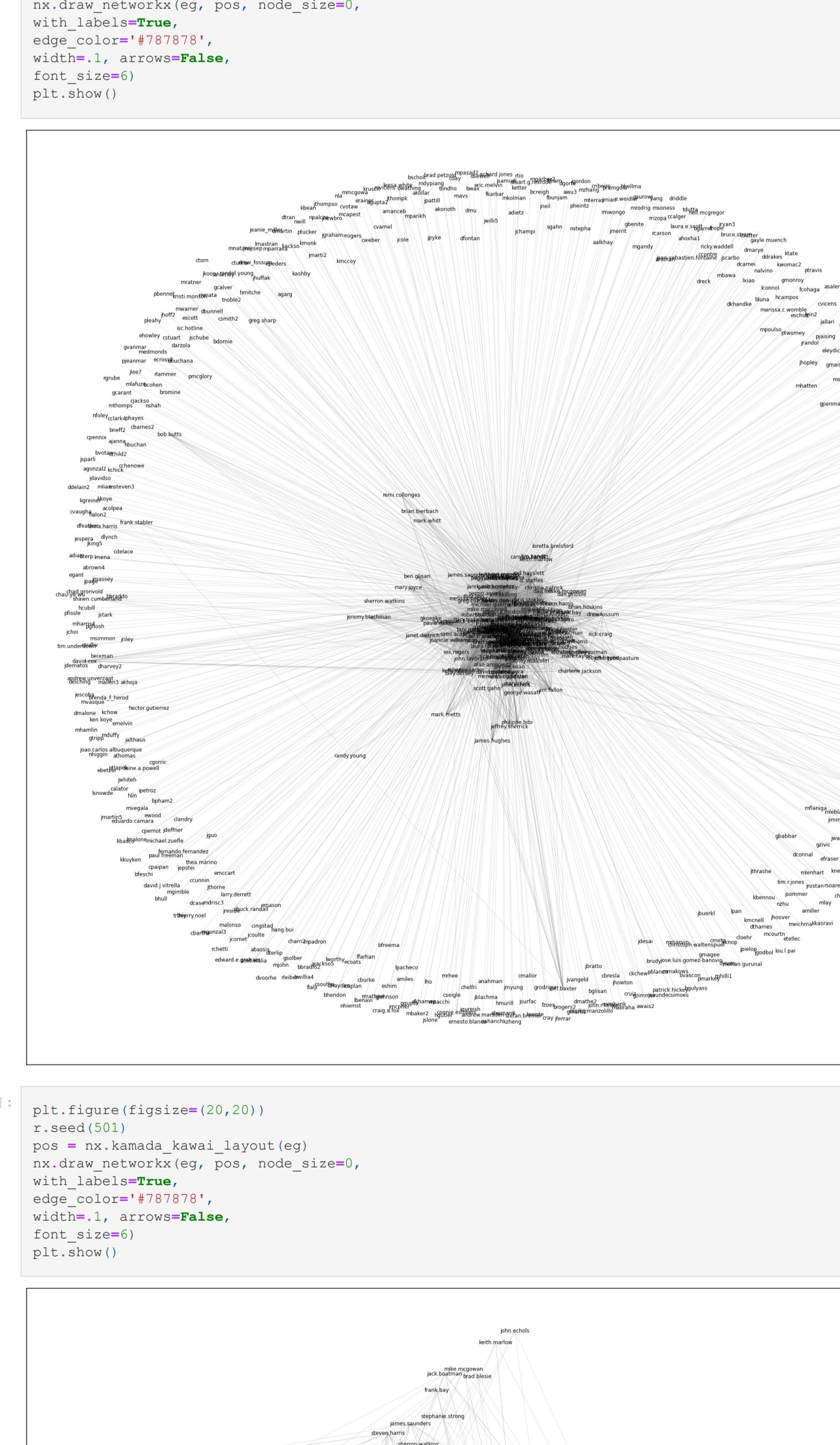
Plotting graph subsets

The email graph is way too large to display the whole thing and get any meaningful information out. However, we can look at subsets of the graph such as the neighbors of a specific node. To visualize it we can use different strategies to layout the nodes. In this case, we will use two different layout strategies: *spring* and *kamada-kawai*. According to [Wikipedia](#), these force directed layout strategies have the characteristic: "...the edges tend to have uniform length (because of the spring forces), and nodes that are not connected by an edge tend to be drawn further apart...".

Use networkx `ego_graph()` method to get a radius=1 neighborhood around `jeff.skilling` and draw the spring graph with a plot that is 20x20 inch so we can see details. Then, draw the same subgraph again using the kamada-kawai layout strategy. Finally, get the neighborhood around `kenneth.lay` and draw kamada-kawai.

In [19]:

```
eg=nx.ego_graph(G, "jeff.skilling", radius=1, center=True, undirected=False, distance=None)
eg_kl=nx.ego_graph(G, "kenneth.lay", radius=1, center=True, undirected=False, distance=None)
```



The figure is a network graph illustrating connections between various actors. The nodes are represented by labels placed near their respective points on the grid. The connections are shown as thin grey lines, with many lines overlapping, suggesting a dense web of relationships.

Nodes visible in the graph:

- robert kilmer
- georgeanne hodges
- brian stalling
- kerry roper
- h. boots
- w. vickers
- drew fossum
- heath schisseler
- eric gadd
- shelley corman
- sheila tweed
- mitch taylor
- rick craig
- michael teraso
- melissa becker

The figure displays a network graph with approximately 40 nodes, each labeled with a unique name. The nodes are arranged in a roughly circular pattern, with some central and peripheral clusters. The edges represent connections between nodes, forming a complex web of links. Notable connections include kevin.greiner and greg.hermans having many outgoing edges to other nodes, while tani.nath has fewer but more direct connections.

A complex network graph showing connections between numerous individuals, likely actors or public figures, based on their names. The nodes are represented by small circles, and the connections are shown as thin gray lines. The names are labels placed near their respective nodes. The graph is highly interconnected, with many nodes having multiple connections to others. Some names are more prominent at the top and bottom, while others are scattered throughout the middle.

