# P R O J E C T

Click on below link to download dataset:

## importing libraries

In [1]:
```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import pycountry
import plotly.express as px
from wordcloud import WordCloud
import warnings; warnings.filterwarnings('ignore')
```

## reading dataset

In [2]:
```python
survey_df = pd.read_csv('sods2022/survey_results_public.csv')
```

In [3]:
```python
schema_df = pd.read_csv('sods2022/survey_results_schema.csv')
```

```
In [4]:    1  survey_df.head()
```

Out[4]:

| | ResponseId | MainBranch | Employment | RemoteWork | CodingActivities | EdLevel | LearnCode | |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | None of these | NaN | NaN | NaN | NaN | NaN | |
| **1** | 2 | I am a developer by profession | Employed, full-time | Fully remote | Hobby;Contribute to open-source projects | NaN | NaN | |
| **2** | 3 | I am not primarily a developer, but I write co... | Employed, full-time | Hybrid (some remote, some in-person) | Hobby | Master's degree (M.A., M.S., M.Eng., MBA, etc.) | Books / Physical media;Friend or family member... | documenta |
| **3** | 4 | I am a developer by profession | Employed, full-time | Fully remote | I don't code outside of work | Bachelor's degree (B.A., B.S., B.Eng., etc.) | Books / Physical media;School (i.e., Universit... | doc |
| **4** | 5 | I am a developer by profession | Employed, full-time | Hybrid (some remote, some in-person) | Hobby | Bachelor's degree (B.A., B.S., B.Eng., etc.) | Other online resources (e.g., videos, blogs, f... | doc |

5 rows × 79 columns

---

## preprocessing

---

we need column *qname* as index of `schema_df` DataFrame

```
In [4]:    1  schema_df.set_index('qname', inplace=True)
```

```
In [5]:    1  schema_df = schema_df.question
```

```
In [ ]:    1
```

only data in *question* column is useful, so we will delete other columns

```
In [ ]:    1
```

After deletion

```
In [ ]:    1  schema_df
```

```
In [9]:  1  schema_df.index
```

```
Out[9]: Index(['S0', 'MetaInfo', 'S1', 'MainBranch', 'Employment', 'RemoteWork',
               'CodingActivities', 'S2', 'EdLevel', 'LearnCode', 'LearnCodeOnline',
               'LearnCodeCoursesCert', 'YearsCode', 'YearsCodePro', 'DevType',
               'OrgSize', 'PurchaseInfluence', 'BuyNewTool', 'Country', 'Currency',
               'CompTotal', 'CompFreq', 'S3', 'Language', 'Database', 'Platform',
               'Webframe', 'MiscTech', 'ToolsTech', 'NEWCollabTools', 'OpSys',
               'VersionControlSystem', 'VCInteraction', 'VCHosting',
               'OfficeStackAsync', 'OfficeStackSync', 'Blockchain', 'S4', 'NEWSOSites',
               'SOVisitFreq', 'SOAccount', 'SOPartFreq', 'SOComm', 'S5', 'Age',
               'Gender', 'Trans', 'Sexuality', 'Ethnicity', 'Accessibility',
               'MentalHealth', 'S6', 'TBranch', 'ICorPM', 'WorkExp', 'Knowledge',
               'Frequency', 'TimeSearching', 'TimeAnswering', 'Onboarding',
               'ProfessionalTech', 'SOTeamsUsage', 'TrueFalse', 'S7', 'SurveyLength',
               'SurveyEase', 'Knowledge_1', 'Knowledge_2', 'Knowledge_3',
               'Knowledge_4', 'Knowledge_5', 'Knowledge_6', 'Knowledge_7',
               'Frequency_1', 'Frequency_2', 'Frequency_3', 'TrueFalse_1',
               'TrueFalse_2', 'TrueFalse_3'],
              dtype='object', name='qname')
```

**plot function**

```python
def custom_plot(series, plot_height=15, plot_width=5,
                y_label_font_size=13.5,
                title = '', title_font_size=15,
                percent_font_size=14,
                color = 'light:#59C1BD'):

    # create figure to display plot
    plt.figure( figsize=(plot_width, plot_height) )

    # to hide square of the plot
    custom_params = {
                    "axes.spines.bottom": False,
                    "axes.spines.right": False,
                    "axes.spines.left"  : False,
                    "axes.spines.top": False
                    }

    sns.set_theme(style="white", rc=custom_params)

    # creating different shades of colors(color palette) of size series leng
    # pal stores rgb values for different color shades
    pal = sns.color_palette(color, len(series)) # light:#5A9

    # argsort return indices of elements according to sorting order..
    # means lowest number will be indexed as 0, and so on
    # rank stores rank of series whr highest count value comes first
    # using this rank to assign color shades to diffrnt bars in plot
    rank = series.argsort().argsort()

    ax = sns.barplot(x = series.values, y=series.index,
                    #palette='PuBuGn_r'
                    #order=series.sort_values('Growth').State,
                     palette=np.array(pal[::])[rank]
                    )

    # to calculate percentage
    s = series.values.sum()

    for rect in ax.patches:
        x_value = rect.get_width()
        y_value = (rect.get_y() + rect.get_height() / 2)
        space = 0

        # calculating percentage and assigning to variable label
        label = "{:.2f}%".format( (100*x_value/s))

        # to display percentage value on bar
        plt.annotate(
            text=label,                       # Use `label` as label
            xy=((x_value/2)-5, y_value),   # Place label at end of the bar,
            xytext=(space, 0),             # Horizontally shift label by `sp
            textcoords="offset points",    # Interpret `xytext` as offset in
            va='center',                   # Vertically center label
            color = 'white',
            #ha='center',
            weight='bold', size=percent_font_size
        )

    plt.title('\n'+title+'\n',
                fontdict=
                {
                    "color": 'black',
                    "weight":'bold',
                    "size":title_font_size
                }
```

```
66                 )
67
68
69     plt.yticks(size=y_label_font_size)#, weight='bold')
70     plt.xticks([], []) # to hide xticks
71
72     f_dict={"color": 'black',"weight":'bold',  "size":15}
73     plt.figtext(.74, .042, "Total Responses: {}".format(s),
74               fontdict = f_dict);
```

## what is your main branch..?

In [11]:
```
1  def MainBranch_ylabel_text_process(s):
2      if s == 'I am not primarily a developer, but I write code sometimes as p
3          return 'Not developer, write\n code as part of work'
4      elif s == 'I used to be a developer by profession, but no longer am':
5          return 'developer by profession\n but no longer'
6      elif s == 'I am a developer by profession':
7          return 'I am a developer\n by profession'
8      elif s == 'I code primarily as a hobby':
9          return 'I code primarily\n as a hobby'
10     else:
11         return s
```

In [12]:
```
1  survey_df['MainBranch'] = survey_df.MainBranch.apply(MainBranch_ylabel_text_
2
3  mb = survey_df.MainBranch.value_counts()
4
5  custom_plot(
6              mb, plot_height=5, y_label_font_size=10, plot_width=12,
7              color = 'light:#000C66'
8          )
```



## How old is the average professional developer..?

In [13]:
```
1  # Age
```

```
1  reorder_list = ['Under 18 years old', '18-24 years old',
2                  '25-34 years old', '35-44 years old',
3                   '45-54 years old', '55-64 years old',
4                   '65 years or older', 'Prefer not to say']
5
6  age_data = survey_df.Age.value_counts().reindex( reorder_list )
7
8  custom_plot(age_data, plot_height=5, color='light:#000C66',
9              title = schema_df.Age, plot_width=10)
```

**What is your age?**

| Age | |
|---|---|
| Under 18 years old | 5.45% |
| 18-24 years old | 23.46% |
| 25-34 years old | 39.62% |
| 35-44 years old | 19.72% |
| 45-54 years old | 7.44% |
| 55-64 years old | 2.7 |
| 65 years or older | |
| Prefer not to say | |

**Total Responses: 70946**

## Employment status of an employee

```
1  # Employment
```

```
1  def colum_expand( s ):
2      d = {}
3
4      for t in s.dropna().values:
5          for i in t.split(';'):
6              if i in d.keys():
7                  d[i] += 1
8              else:
9                  d[i] = 1
10
11     return pd.Series(d)
12
```

```
1  emp = colum_expand(survey_df.Employment)
2
3  custom_plot(emp, plot_height=7, color='light:#000C66',
4          title=schema_df.Employment, plot_width=9)
```

**Which of the following best describes your current employment status?**

| | |
|---|---|
| Employed, full-time | 58.32% |
| Student, full-time | 12.96% |
| Student, part-time | 4. |
| Not employed, but looking for work | 4. |
| Independent contractor, freelancer, or self-employed | 12.71% |
| Employed, part-time | 4.9 |
| Not employed, and not looking for work | |
| Retired | |
| I prefer not to say | |

**Total Responses: 84360**

In [ ]:

```
1
```

# mode of working of employee(remote/hybrid)

In [ ]:

```
1  # RemoteWork
```

In [27]:

```
1  remote_work = survey_df.RemoteWork.value_counts()
2
3  custom_plot(remote_work, plot_height=2.1, plot_width=9
4          , color='dark:#000C66', title=schema_df.RemoteWork)
```

**Which best describes your current work situation?**

| RemoteWork | |
|---|---|
| Fully remote | 42.98% |
| Hybrid (some remote, some in-person) | 42.44% |
| Full in-person | 14.58% |

**Total Responses: 58958**

# how many of you write code outside of your work

```
In [30]:    1  coding_act = colum_expand(survey_df.CodingActivities)
            2
            3  custom_plot(coding_act, plot_height=5, plot_width=9,
            4          color='light:#000C66', title=schema_df.CodingActivities)
```

**Which of the following best describes the code you write outside of work? Select all that apply.**

| | |
|---|---|
| Hobby | 43.77% |
| Contribute to open-source projects | 15.68% |
| I don't code outside of work | 7.46% |
| Bootstrapping a business | 8.57% |
| Freelance/contract work | 13.57% |
| Other (please specify): | 2 |
| School or academic work | 8.73% |

**Total Responses: 98057**

```
In [ ]:    1
```

# What is your highest level of formal education..?

```
In [67]:    1  edu = survey_df.EdLevel.value_counts()
            2
            3  custom_plot(edu, plot_height=9, plot_width=7, color='light:#000C66',
            4          title=schema_df.EdLevel)
```

**Which of the following best describes the highest level of formal education that you've completed? ***

| | |
|---|---|
| Bachelor's degree (B.A., B.S., B.Eng., etc.) | 42.30% |
| Master's degree (M.A., M.S., M.Eng., MBA, etc.) | 21.64% |
| Some college/university study without earning a degree | 13.03% |
| Secondary school (e.g. American high school, German Realschule or Gymnasium, etc.) | 11.04% |
| Associate degree (A.A., A.S., etc.) | 3. |
| Other doctoral degree (Ph.D., Ed.D., etc.) | 3. |
| Primary/elementary school | 2 |
| Something else | 1 |
| Professional degree (JD, MD, etc.) | |

EdLevel

**Total Responses: 71571**

```
In [ ]:    1
```

# How did you learn to code

```
1  learn_code_data = colum_expand(survey_df.LearnCode)
2
3
4
5  custom_plot(learn_code_data, plot_height=9, plot_width=10,
6          color='light:#000C66', title=schema_df.LearnCode,
7          y_label_font_size=20)
8
```

**How did you learn to code? Select all that apply.**

| | |
|---|---|
| Books / Physical media | 16.53% |
| Friend or family member | 4.23% |
| Other online resources (e.g., videos, blogs, forum) | 21.52% |
| School (i.e., University, College, etc) | 18.87% |
| On the job training | 12.09% |
| Online Courses or Certification | 14.15% |
| Coding Bootcamp | 3.28% |
| Colleague | 5.59% |
| Other (please specify): | 1. |
| Hackathons (virtual or in-person) | 2.2 |

**Total Responses: 235891**

```
1
```

# What online resources do you use to learn to code?

```
In [32]:  1  learn_code_online = colum_expand(survey_df.LearnCodeOnline)
          2
          3
          4  custom_plot(learn_code_online, plot_height=14, plot_width=11,
          5             color='light:#000C66', title=schema_df.LearnCodeOnline)
          6
```

**What online resources do you use to learn to code? Select all that apply.**

| Resource | Percentage |
|---|---|
| Technical documentation | 13.45% |
| Blogs | 11.50% |
| Programming Games | 2.03% |
| Written Tutorials | 8.86% |
| Stack Overflow | 13.15% |
| Online books | 6.70% |
| Video-based Online Courses | 7.85% |
| Online challenges (e.g., daily or weekly coding challenges) | 3.83% |
| Online forum | 6.16% |
| How-to videos | 9.14% |
| Written-based Online Courses | 5.25% |
| Auditory material (e.g., podcasts) | 1.1 |
| Interactive tutorial | 4.00% |
| Coding sessions (live or recorded) | 4.40% |
| Other (Please specify): | |
| Certification videos | 2.27% |

**Total Responses: 332107**

```
In [ ]:   1
```

## What online courses or certifications do you use to learn to code?

```
In [36]:  1  learn_code_cert = colum_expand(survey_df.LearnCodeCoursesCert)
          2
          3  custom_plot(learn_code_cert, plot_height=5, plot_width=10,
          4            color='light:#000C66', title=schema_df.LearnCodeCoursesCert)
```
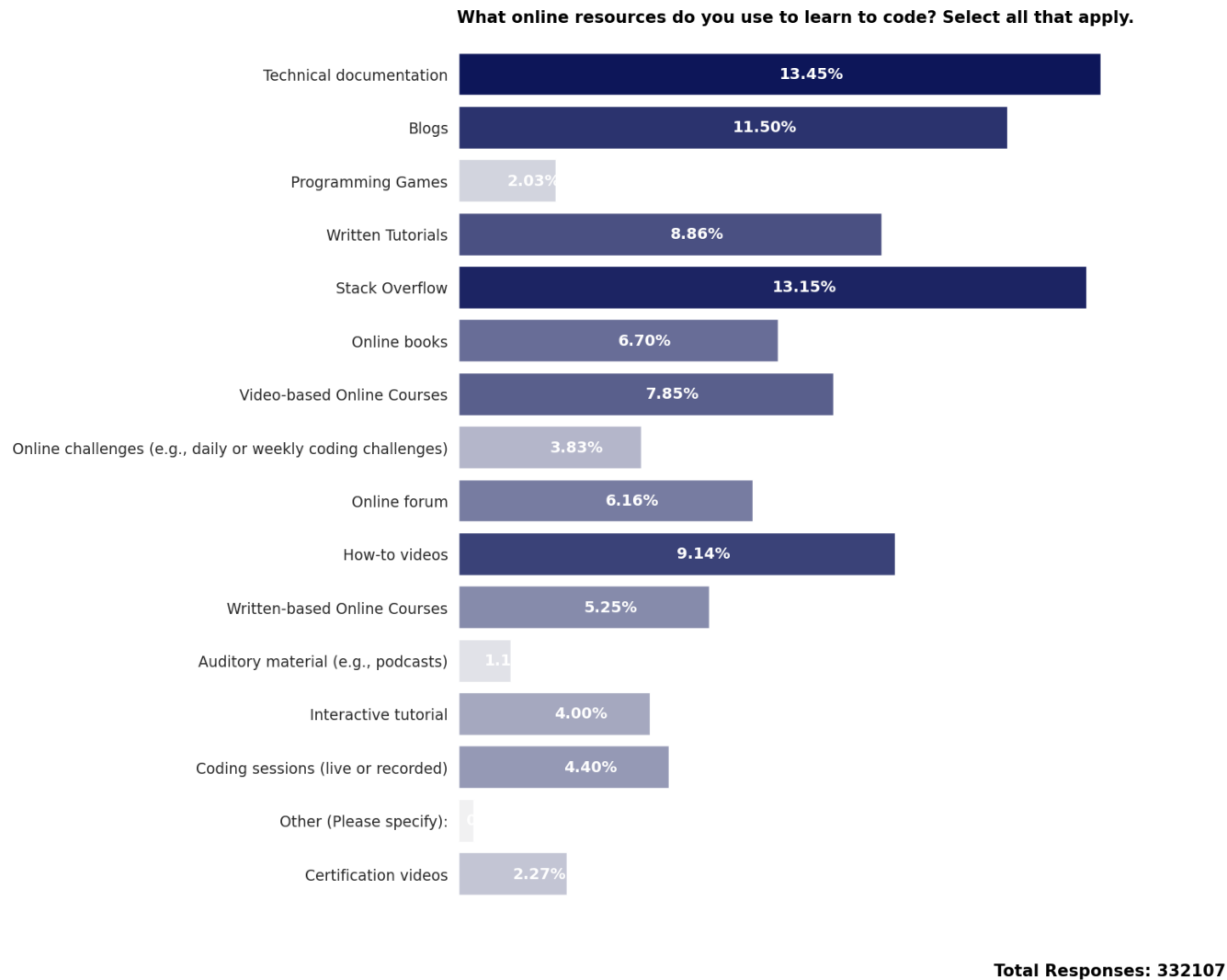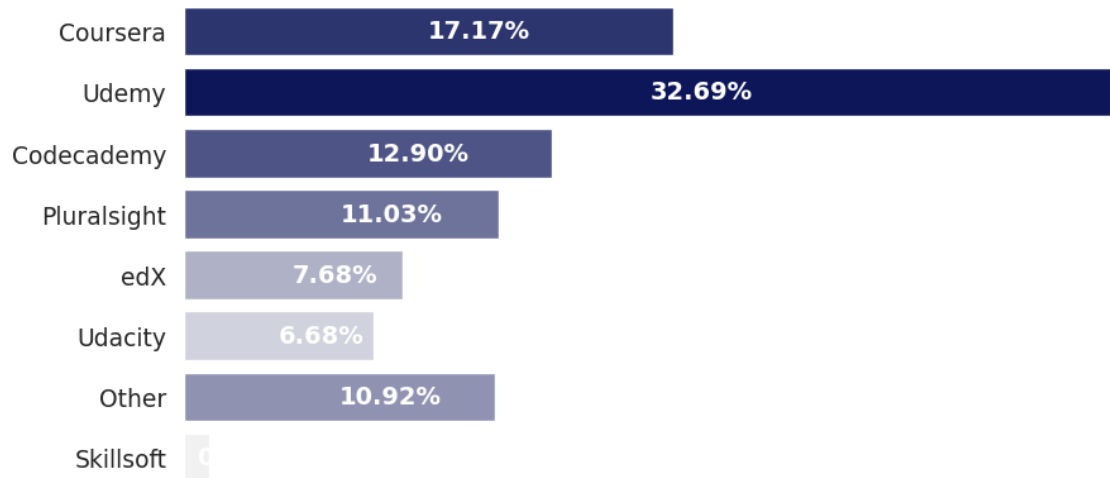
**What online courses or certifications do you use to learn to code? Select all that apply.**

Coursera — 17.17%
Udemy — 32.69%
Codecademy — 12.90%
Pluralsight — 11.03%
edX — 7.68%
Udacity — 6.68%
Other — 10.92%
Skillsoft

**Total Responses: 59773**

```
In [ ]:  1
```

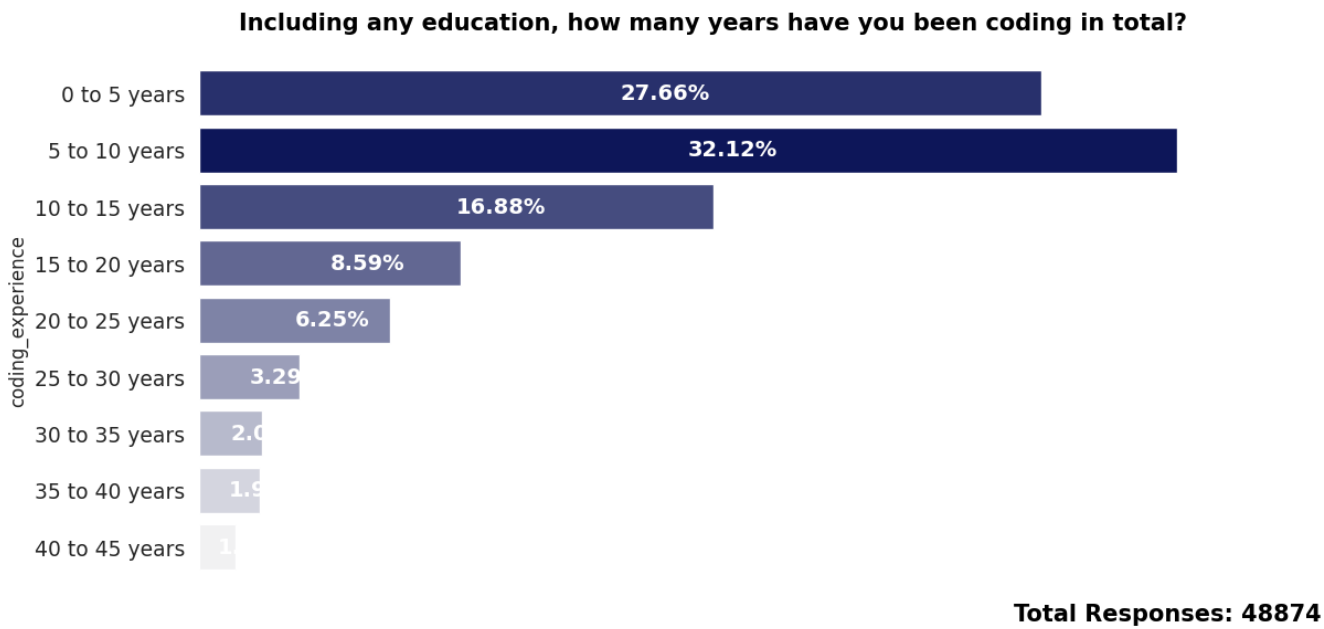## how many years have you been coding in total (Including education)

```
In [37]:  1  def make_groups(s):
          2      try:
          3          s = int(s)
          4          if s > 0 and s < 5:
          5              return '0 to 5 years'
          6          if s > 5 and s < 10:
          7              return '5 to 10 years'
          8          if s > 10 and s < 15:
          9              return '10 to 15 years'
         10          if s > 15 and s < 20:
         11              return '15 to 20 years'
         12          if s > 20 and s < 25:
         13              return '20 to 25 years'
         14          if s > 25 and s < 30:
         15              return '25 to 30 years'
         16          if s > 30 and s < 35:
         17              return '30 to 35 years'
         18          if s > 35 and s < 40:
         19              return '35 to 40 years'
         20          if s > 40 and s < 45:
         21              return '40 to 45 years'
         22          if s > 45 and s < 50:
         23              return '45 to 50 years'
         24      except (TypeError, ValueError):
         25          pass
         26
```

```
In [38]:  1  survey_df['coding_experience'] = survey_df.YearsCode.apply(make_groups)
```

```
1
2  reorder_list = ['0 to 5 years', '5 to 10 years', '10 to 15 years',
3                  '15 to 20 years', '20 to 25 years', '25 to 30 years',
4                  '30 to 35 years','35 to 40 years', '40 to 45 years']
5
6  ce = survey_df.coding_experience.value_counts().reindex(reorder_list)
7
8  custom_plot(ce, plot_height=6, plot_width=12,
9              title=schema_df.YearsCode, color='light:#000C66')
```

**Including any education, how many years have you been coding in total?**

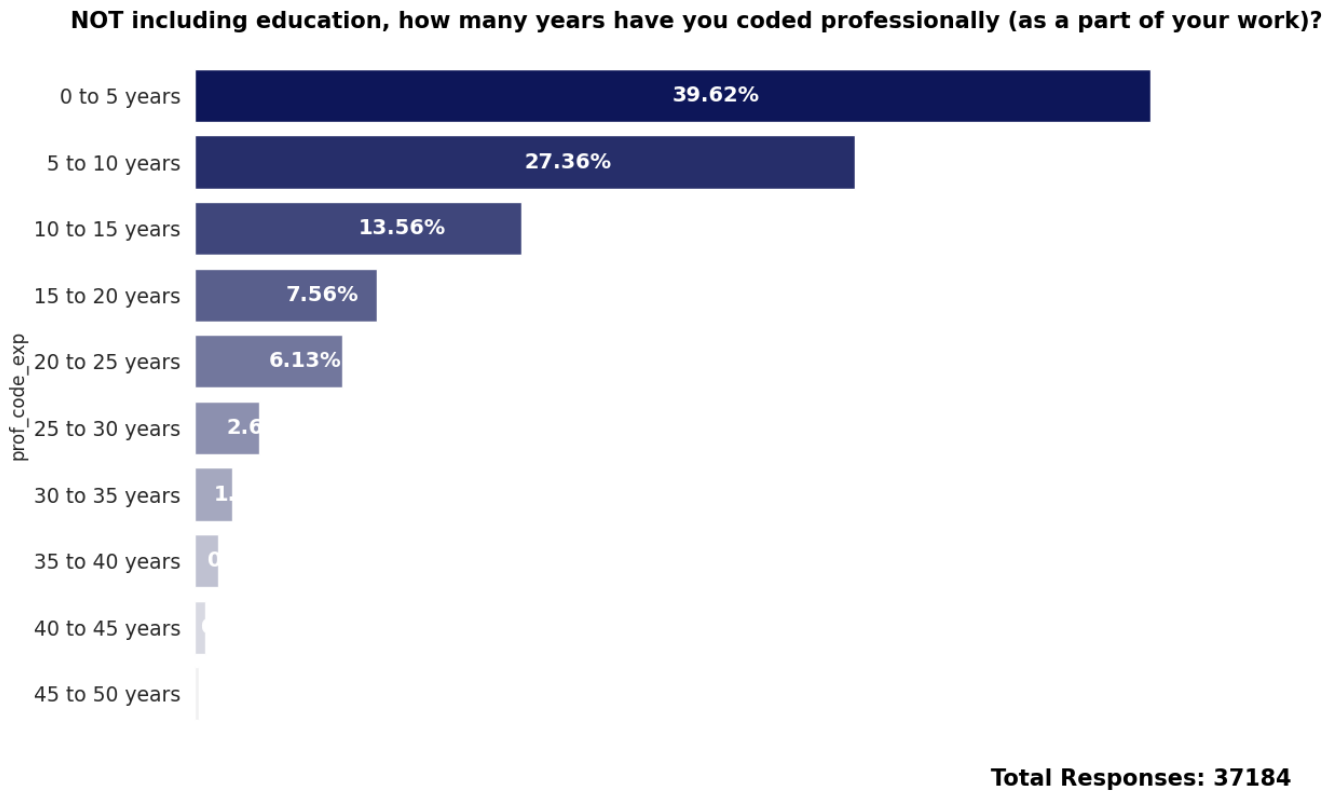| | |
|---|---|
| 0 to 5 years | 27.66% |
| 5 to 10 years | 32.12% |
| 10 to 15 years | 16.88% |
| 15 to 20 years | 8.59% |
| 20 to 25 years | 6.25% |
| 25 to 30 years | 3.29 |
| 30 to 35 years | 2.0 |
| 35 to 40 years | 1.9 |
| 40 to 45 years | 1 |

**Total Responses: 48874**

## how many years have you been coding in total ( not Including education)

In [ ]:
```
1  # YearsCodePro
```

In [46]:
```
1  survey_df['prof_code_exp'] = survey_df.YearsCodePro.dropna().apply(make_grou
```

```
In [49]:    1  pce = survey_df.prof_code_exp.value_counts()
            2
            3  custom_plot(pce, plot_height=8,plot_width=12,
            4          title=schema_df.YearsCodePro, color='light:#000C66')
```

**NOT including education, how many years have you coded professionally (as a part of your work)?**



Total Responses: 37184

## what kind of developer you are..?

```
In [ ]:    1  # DevType
```

```
In [50]:    1  schema_df.DevType
```

```
Out[50]: 'Which of the following describes your current job? Please select all that appl
         y.'
```

```
In [51]:    1  survey_df.DevType
```

```
Out[51]: 0                                                       NaN
         1                                                       NaN
         2          Data scientist or machine learning specialist;...
         3                                       Developer, full-stack
         4          Developer, front-end;Developer, full-stack;Dev...
                                    ...
         73263                                  Developer, back-end
         73264        Data scientist or machine learning specialist
         73265     Developer, full-stack;Developer, desktop or en...
         73266     Developer, front-end;Developer, desktop or ent...
         73267     Developer, front-end;Engineer, data;Engineer, ...
         Name: DevType, Length: 73268, dtype: object
```

```
In [53]:   1  dev_type = colum_expand(survey_df.DevType)
           2
           3  dev_type = dev_type.sort_values(ascending=False)
           4
           5  custom_plot(dev_type, color='light:#000C66', plot_height=28,
           6              plot_width=14, y_label_font_size=22)
```

| | |
|---|---|
| Developer, full-stack | 17.42% |
| Developer, back-end | 16.14% |
| Developer, front-end | 9.66% |
| Developer, desktop or enterprise applications | 5.79% |
| Developer, mobile | 4.63% |
| DevOps specialist | 3.74% |
| Student | 3.40% |
| Cloud infrastructure engineer | 3.21% |
| Database administrator | 2.99% |
| System administrator | 2.98% |
| Developer, embedded applications or devices | 2.38% |
| Project manager | 2.36% |
| Designer | 2.28% |
| Engineer, data | 2.18% |
| Engineering manager | 2.17% |
| Data scientist or machine learning specialist | 2.08% |
| Data or business analyst | 1.94% |
| Developer, QA or test | 1.88% |
| Academic researcher | 1.64 |
| Other (please specify): | 1.59 |
| Product manager | 1.53 |
| Educator | 1.2 |
| Engineer, site reliability | 1.1 |
| Security professional | 1.1 |
| Developer, game or graphics | 1.1 |
| Senior Executive (C-Suite, VP, etc.) | 1.1 |
| Scientist | 1.0 |
| Blockchain | 0. |
| Marketing or sales professional | |

**Total Responses: 164790**

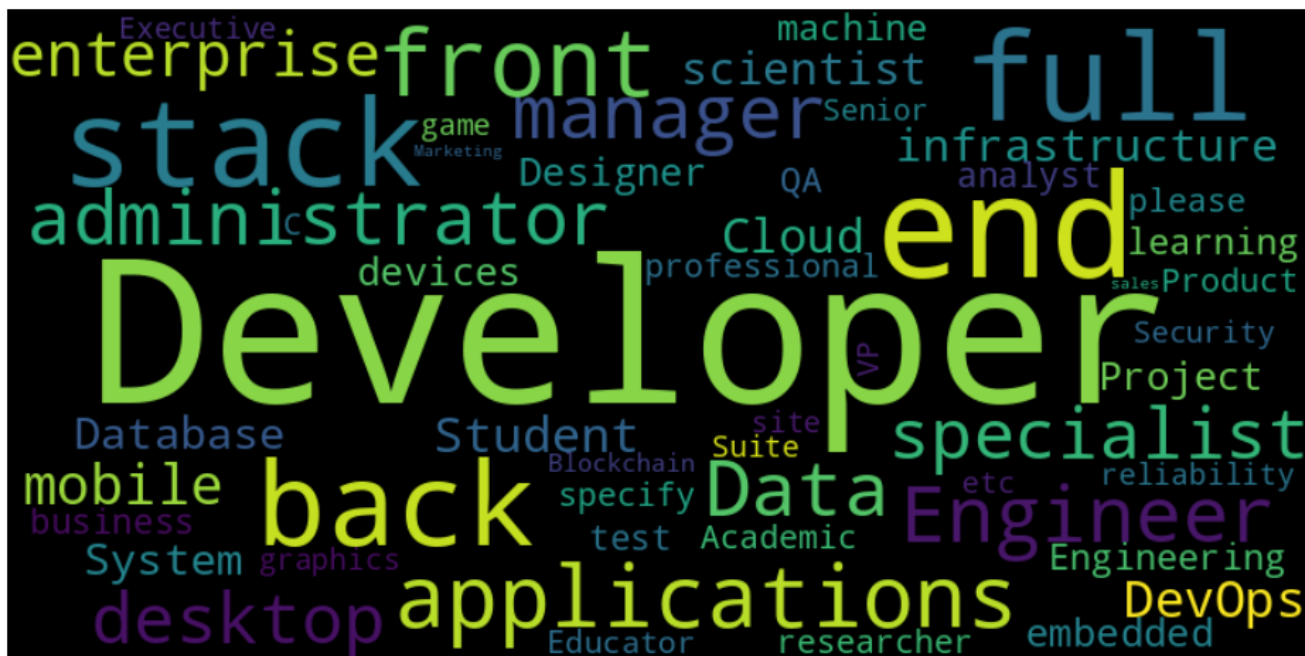# WordCloud for developer type

```
In [58]:   1  words = ' '.join((job for job in survey_df.DevType.dropna().str.replace(';',
```

```
In [ ]:  1
```

```
In [61]:  1  # Generate word cloud
          2  wordcloud = WordCloud(collocation_threshold=int(1e6), width=800, height=400,
          3                        background_color='black').generate(words)
          4
          5  # Plot the word cloud
          6  plt.figure(figsize=(12, 6))
          7  plt.imshow(wordcloud, interpolation='bilinear')
          8  plt.axis('off')
          9  plt.show()
```



```
In [ ]:  1
```

## what is organization size of the developer...?

```
In [ ]:  1  # OrgSize
```
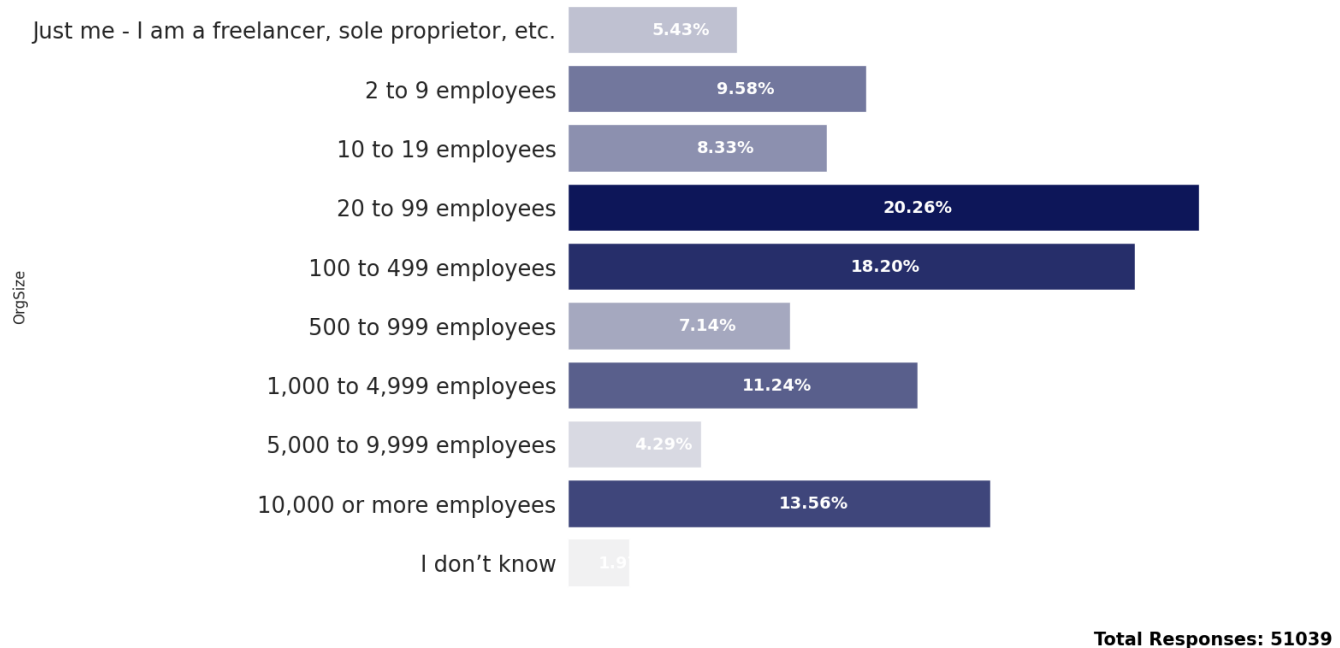
```
In [62]:  1  schema_df.OrgSize
```

```
Out[62]: 'Approximately how many people are employed by the company or organization you
         currently work for? '
```

```
In [65]:  1  survey_df.OrgSize.value_counts()
```

```
Out[65]: OrgSize
         20 to 99 employees                             10343
         100 to 499 employees                            9289
         10,000 or more employees                        6922
         1,000 to 4,999 employees                        5736
         2 to 9 employees                                4887
         10 to 19 employees                              4251
         500 to 999 employees                            3645
         Just me - I am a freelancer, sole proprietor, etc.  2771
         5,000 to 9,999 employees                        2189
         I don't know                                    1006
         Name: count, dtype: int64
```

```
In [66]:  1  reorder_list = [
          2      "Just me - I am a freelancer, sole proprietor, etc.",
          3      "2 to 9 employees", "10 to 19 employees", "20 to 99 employees",
          4      "100 to 499 employees", "500 to 999 employees",
          5      "1,000 to 4,999 employees", "5,000 to 9,999 employees",
          6      "10,000 or more employees", "I don't know"
          7  ]
          8
          9  org_size = survey_df.OrgSize.value_counts().reindex(reorder_list)
         10
         11  custom_plot(org_size, plot_height=9, plot_width=10,
         12              color = 'light:#000C66',
         13              y_label_font_size=18.5)
```



Total Responses: 51039

## donut plot function using pie plot
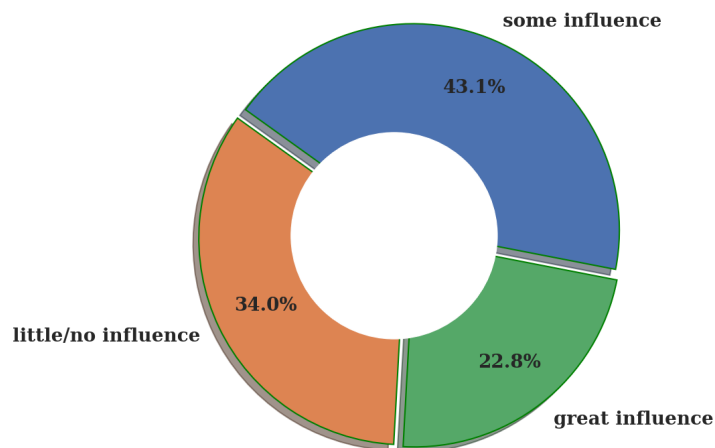
```
In [104]:  1  def plot_pie(data , title='', distance_btwn_pieces=0.09, startangle=-11):
           2
           3      explode = (distance_btwn_pieces,) * len(data)
           4      plt.figure(figsize=(14,10))
           5
           6
           7      plt.pie( data, explode=explode,  labels=data.index, pctdistance=0.75,
           8              colors = ['red', 'blue', 'yellow','pink','blue'],
           9              wedgeprops={'linewidth': 1.5, 'edgecolor' : "green" },
          10              textprops={"weight":'bold', "size":20, 'family':'serif'},
          11              autopct='%1.1f%%',  startangle=startangle, shadow=True,
          12              )
          13
          14      #plt.setp(pcts, color='black')
          15      hfont = {'fontname': 'serif', 'weight': 'bold'}
          16      plt.title(title, size=25, **hfont)
          17
          18      centre_circle = plt.Circle((-0.08,0), 0.5, fc='white')
          19      fig = plt.gcf().gca().add_artist(centre_circle)
          20      ;
          21
          22
```

## What level of influence developer, have over new technology purchases at your organization?

```
In [ ]:    1  # PurchaseInfluence
```

```
In [76]:   1  def shorten_names(s):
           2      if s == 'I have some influence':
           3          return 'some influence'
           4      elif s == 'I have little or no influence':
           5          return 'little/no influence'
           6      elif s == 'I have a great deal of influence':
           7          return 'great influence'
           8
           9  tech_influence = survey_df.PurchaseInfluence.apply(shorten_names)
          10  tech_influence = tech_influence.value_counts()
          11
          12  plot_pie(data=tech_influence,
          13          title = schema_df.PurchaseInfluence,
          14          distance_btwn_pieces=0.03
          15          )
```

**What level of influence do you, personally, have over new technology purchases at your organization?**

some influence

43.1%

34.0%

little/no influence

22.8%

great influence

## Where do developer live?

```
In [79]:   1  country = survey_df.Country.value_counts()[:40]
           2
           3  custom_plot(country, y_label_font_size=15, plot_height=35,
           4              title= schema_df.Country.split('<')[0],
           5              color = 'light:#000C66')
```

## Where do you live?

| Country | Percentage |
|---|---|
| United States of America | 21.70% |
| India | 10.64% |
| Germany | 8.65% |
| United Kingdom of Great Britain and Northern Ireland | 6.72% |
| Canada | 3.9 |
| France | 3.7 |
| Brazil | 3.3 |
| Poland | 2. |
| Netherlands | 2. |
| Spain | 2. |
| Italy | 2. |
| Australia | 2. |
| Russian Federation | 1. |
| Turkey | 1 |
| Sweden | 1 |
| Switzerland | 1 |
| Austria | 1 |
| Israel | 1 |
| Iran, Islamic Republic of... | 1 |
| Pakistan | 1 |
| Czech Republic | 1 |
| China | 1 |
| Belgium | 1 |
| Bangladesh | 0 |
| Ukraine | 0 |
| Romania | 0 |
| Mexico | 0 |
| Portugal | 0 |
| Greece | 0 |

Denmark

Indonesia

Argentina

Nigeria

South Africa

Norway

Finland

Hungary

New Zealand

Egypt
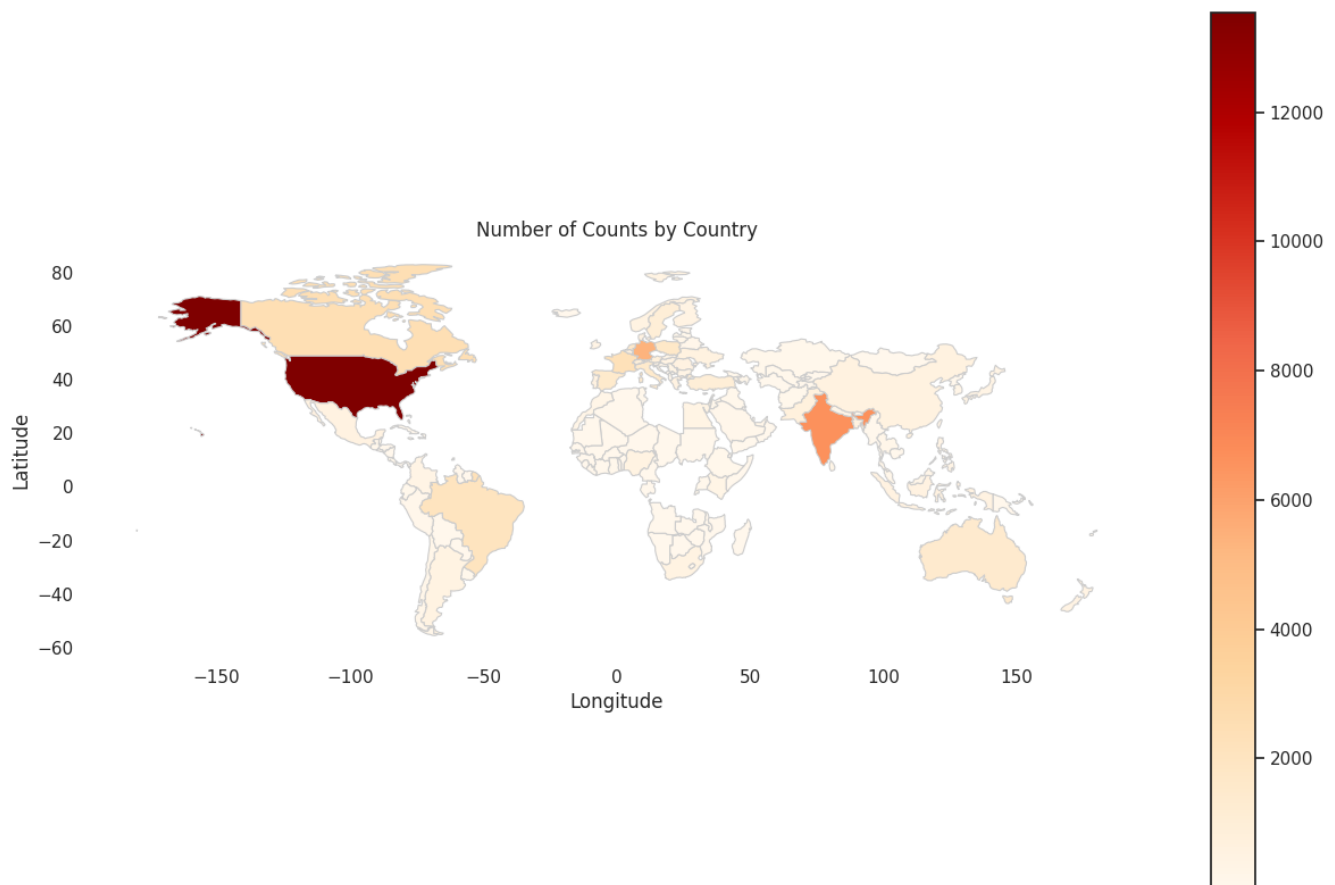
Philippines

**Total Responses: 62397**

# Map plot of country with developer count

```
In [119]:   1  d = survey_df.Country.value_counts().reset_index()
```

```
In [107]:   1  d.to_csv('country_for_map.csv')
```

```
In [84]:
 1  import geopandas as gpd
 2
 3  data = pd.read_csv("country_for_map.csv")
 4
 5  # Load the world map
 6  world = gpd.read_file(gpd.datasets.get_path('naturalearth_lowres'))
 7
 8  # Merge the world map with the data DataFrame
 9  world = world.merge(data, how='left', left_on='name', right_on='Country')
10
11  # Plot the map
12  fig, ax = plt.subplots(1, 1, figsize=(15, 10))
13  world.plot(column='count', cmap='OrRd', linewidth=0.8, ax=ax, edgecolor='0.8
14
15
16  plt.title('Number of Counts by Country')
17  plt.xlabel('Longitude')
18  plt.ylabel('Latitude')
19
20  plt.show()
```



### Which currency does developer use day-to-day?

```
In [ ]:
 1  # Currency
```

```
In [85]:
 1  schema_df.Currency
```

```
Out[85]: "Which currency do you use day-to-day? If your answer is complicated, please pi
         ck the one you're most comfortable estimating in. *"
```
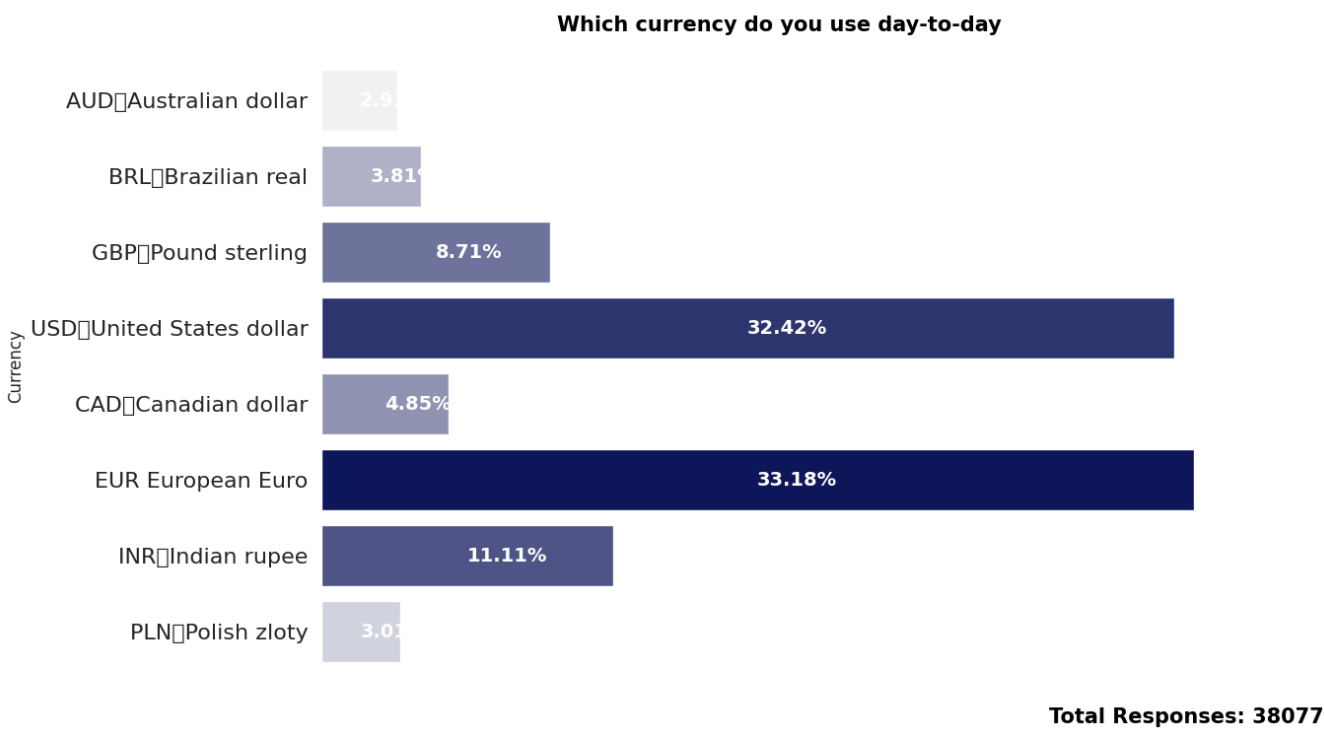
```
In [ ]:
 1
```

```
In [88]:  1  survey_df.Currency.value_counts()
```

```
Out[88]:  Currency
          EUR European Euro                  12634
          USD\tUnited States dollar          12346
          INR\tIndian rupee                   4229
          GBP\tPound sterling                 3318
          CAD\tCanadian dollar                1847
                                              ...
          BND\tBrunei dollar                     1
          PGK\tPapua New Guinean kina             1
          SHP\tSaint Helena pound                 1
          GIP\tGibraltar pound                    1
          TOP\tTongan pa'anga                     1
          Name: count, Length: 142, dtype: int64
```

```
In [97]:  1  currency = survey_df.Currency.value_counts()[:8]
          2  currency = currency.sample(len(currency))
          3
          4  custom_plot(currency, plot_height=8, plot_width=12,
          5              color = 'light:#000C66',
          6              title=schema_df.Currency.split('?')[0],
          7              y_label_font_size=16)
```

**Which currency do you use day-to-day**



Total Responses: 38077

# VersionControlSystem

```
In [ ]:  1  # VersionControlSystem
```

```
In [98]:  1  schema_df.VersionControlSystem
```

```
Out[98]:  'What are the primary <b>version control systems</b> you use? Select all that a
          pply.'
```

```
1  vcs = colum_expand(survey_df.VersionControlSystem)
2
3  custom_plot(vcs, plot_height=6, plot_width=13, color = 'light:#000C66',
4            y_label_font_size=15, title=schema_df.VersionControlSystem)
```

**What are the primary <b>version control systems</b> you use? Select all that apply.**

| | |
|---|---|
| Git | 87.43% |
| Other (please specify): | 2 |
| Mercurial | |
| SVN | 4.8 |
| I don't use one | 4. |

**Total Responses: 76641**

```
1  plot_pie(vcs, startangle=25,
2           distance_btwn_pieces=0.1)
```



Git 87.4%

I don't use one 4.0%

SVN 4.8%

Mercurial 1.1%

Other (please specify): 2.7%

# what is your gender..?

```
1  # Gender
```

```
In [108]:    1  colum_expand(survey_df.Gender)
```

```
Out[108]:  Man                                                        65097
           Or, in your own words:                                       521
           Woman                                                       3662
           Non-binary, genderqueer, or gender non-conforming          1186
           Prefer not to say                                           1172
           dtype: int64
```

```
In [109]:    1  gender = colum_expand(survey_df.Gender)
             2
             3  gender.rename( lambda x: x.split(',')[0], inplace=True )
             4
             5  plot_pie(gender,distance_btwn_pieces=0.09, startangle=15,
             6          title=schema_df.Gender)
```

**Which of the following describe you, if any? Please check all that apply.**



# ethincity of developer

```
In [ ]:      1  # Ethnicity
```

```
In [110]:    1  schema_df.Ethnicity
```

```
Out[110]:  'Which of the following describe you, if any? Please check all that apply.'
```

```
In [113]:  1  colum_expand(survey_df.Ethnicity)
```

```
Out[113]:  White                                                       27360
           Or, in your own words:                                       1524
           Indian                                                       6739
           European                                                    25877
           North American                                               2331
           Middle Eastern                                               2850
           Ethnoreligious group                                          348
           Prefer not to say                                            1732
           African                                                      2294
           Asian                                                        6586
           East Asian                                                   1214
           Black                                                        1028
           Caribbean                                                     460
           Southeast Asian                                              1618
           Central American                                              416
           North African                                                 611
           Hispanic or Latino/a                                         3967
           South American                                               2624
           South Asian                                                  1797
           I don't know                                                  701
           Multiracial                                                  1222
           Biracial                                                      798
           Indigenous (such as Native American or Indigenous Australian)  330
           Pacific Islander                                              147
           Central Asian                                                 397
           dtype: int64
```
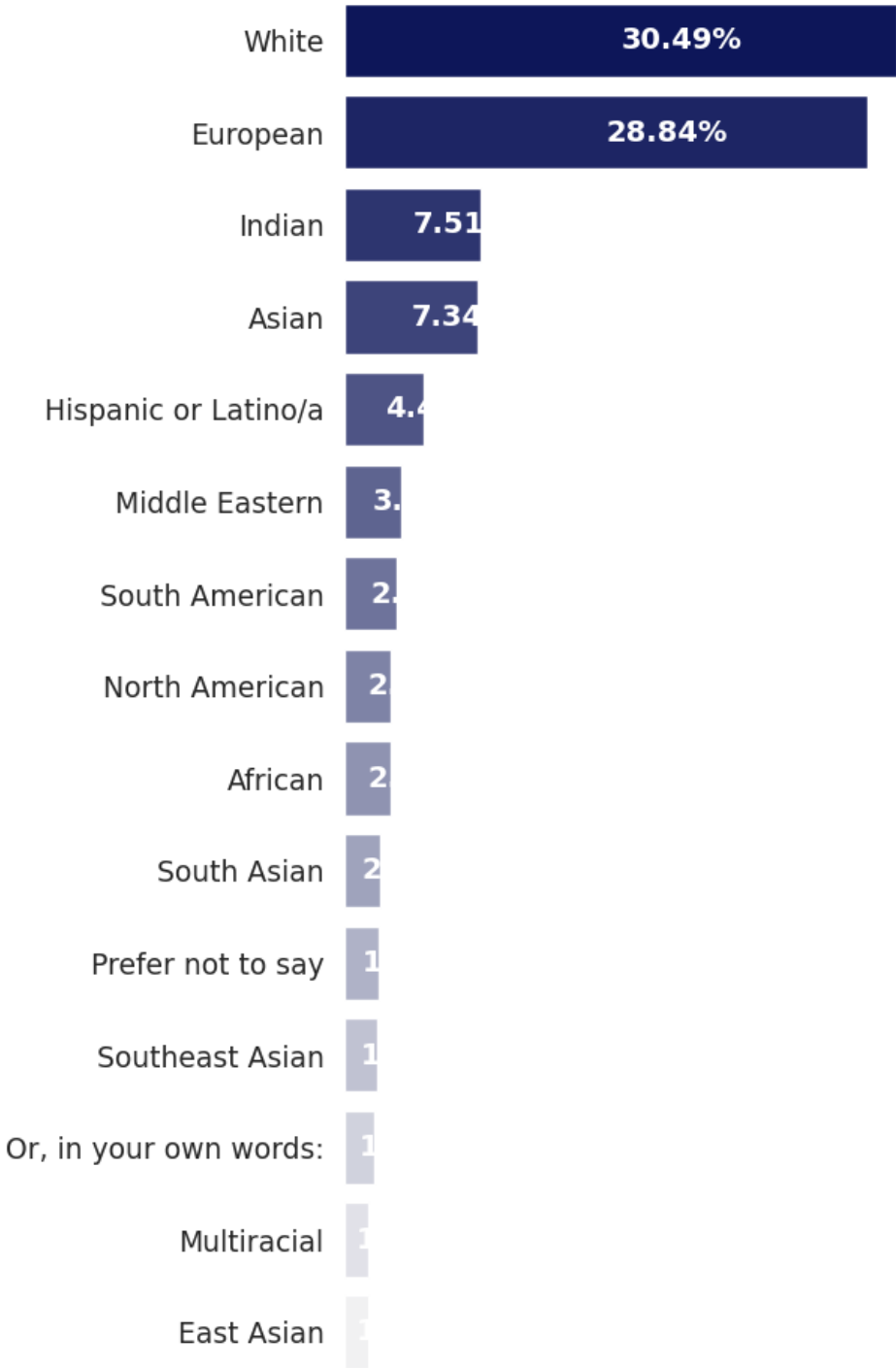
```
1  ethincity = colum_expand(survey_df.Ethnicity).nlargest(15)
2
3  custom_plot(ethincity, plot_height=12, title=schema_df.Ethnicity, color = 'l
```
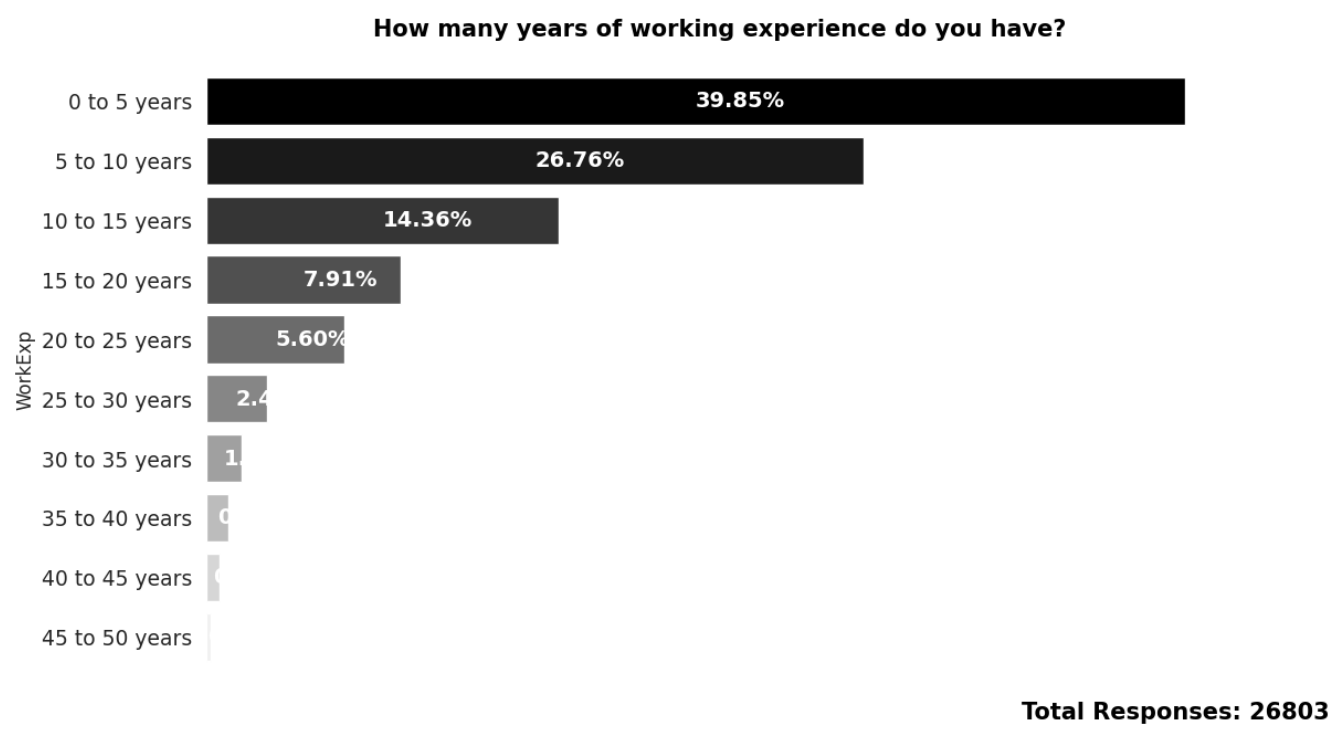
**Which of the following describe you, if any? Please check all that apply.**

| | |
|---|---|
| White | 30.49% |
| European | 28.84% |
| Indian | 7.51 |
| Asian | 7.34 |
| Hispanic or Latino/a | 4.4 |
| Middle Eastern | 3. |
| South American | 2. |
| North American | 2 |
| African | 2 |
| South Asian | 2 |
| Prefer not to say | 1 |
| Southeast Asian | 1 |
| Or, in your own words: | 1 |
| Multiracial | |
| East Asian | |

**Total Responses: 89735**

```
1  # WorkExp
```
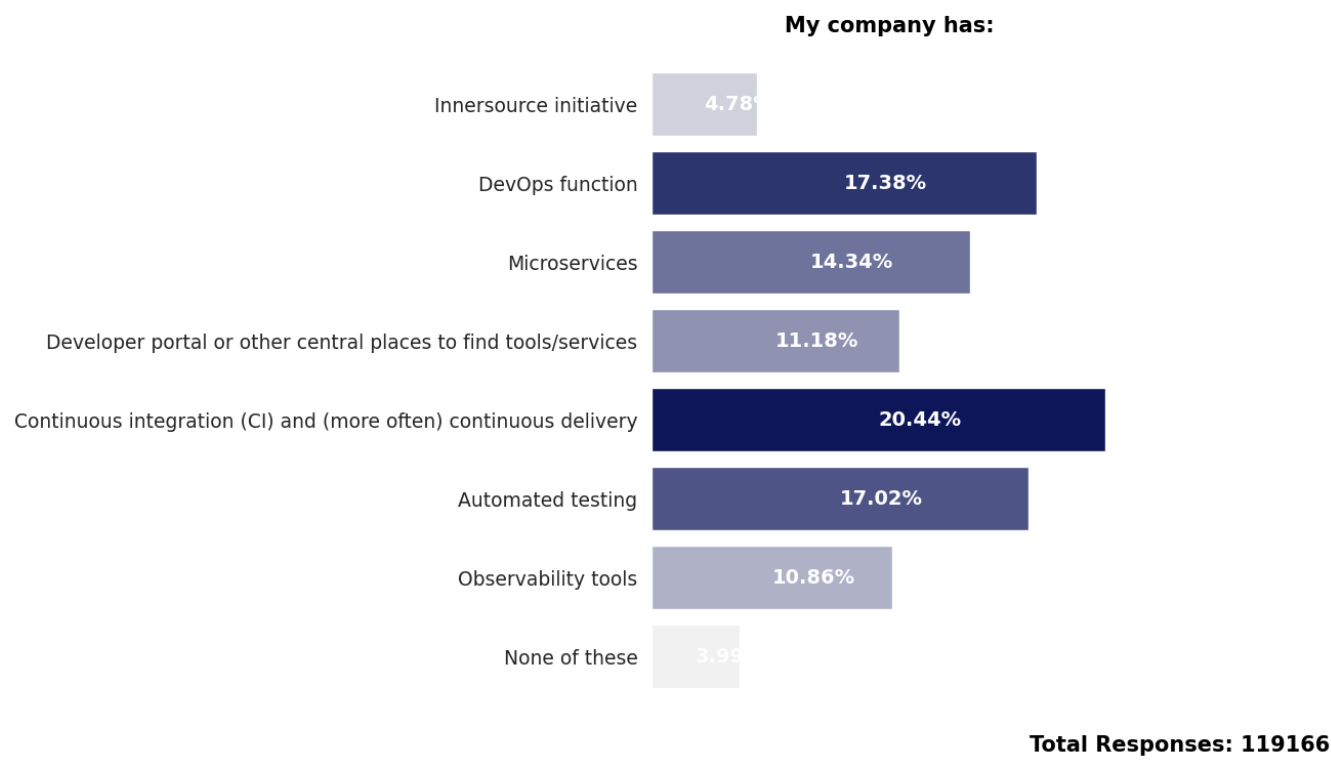
# How many years of working experience do you have?

```
1  work_exp = survey_df.WorkExp.apply(make_groups).value_counts()
2
3
4  custom_plot(work_exp, plot_height=7, plot_width=12, title=schema_df.WorkExp,
```

**How many years of working experience do you have?**

| WorkExp | |
|---|---|
| 0 to 5 years | 39.85% |
| 5 to 10 years | 26.76% |
| 10 to 15 years | 14.36% |
| 15 to 20 years | 7.91% |
| 20 to 25 years | 5.60% |
| 25 to 30 years | 2.4 |
| 30 to 35 years | 1. |
| 35 to 40 years | |
| 40 to 45 years | |
| 45 to 50 years | |

**Total Responses: 26803**

# which technologies does your company have?

```
1  # ProfessionalTech
```

```
1  tech = colum_expand(survey_df.ProfessionalTech)
2
3  custom_plot(tech, plot_height=8, plot_width=6, color = 'light:#000C66',
4              title=schema_df.ProfessionalTech)
```
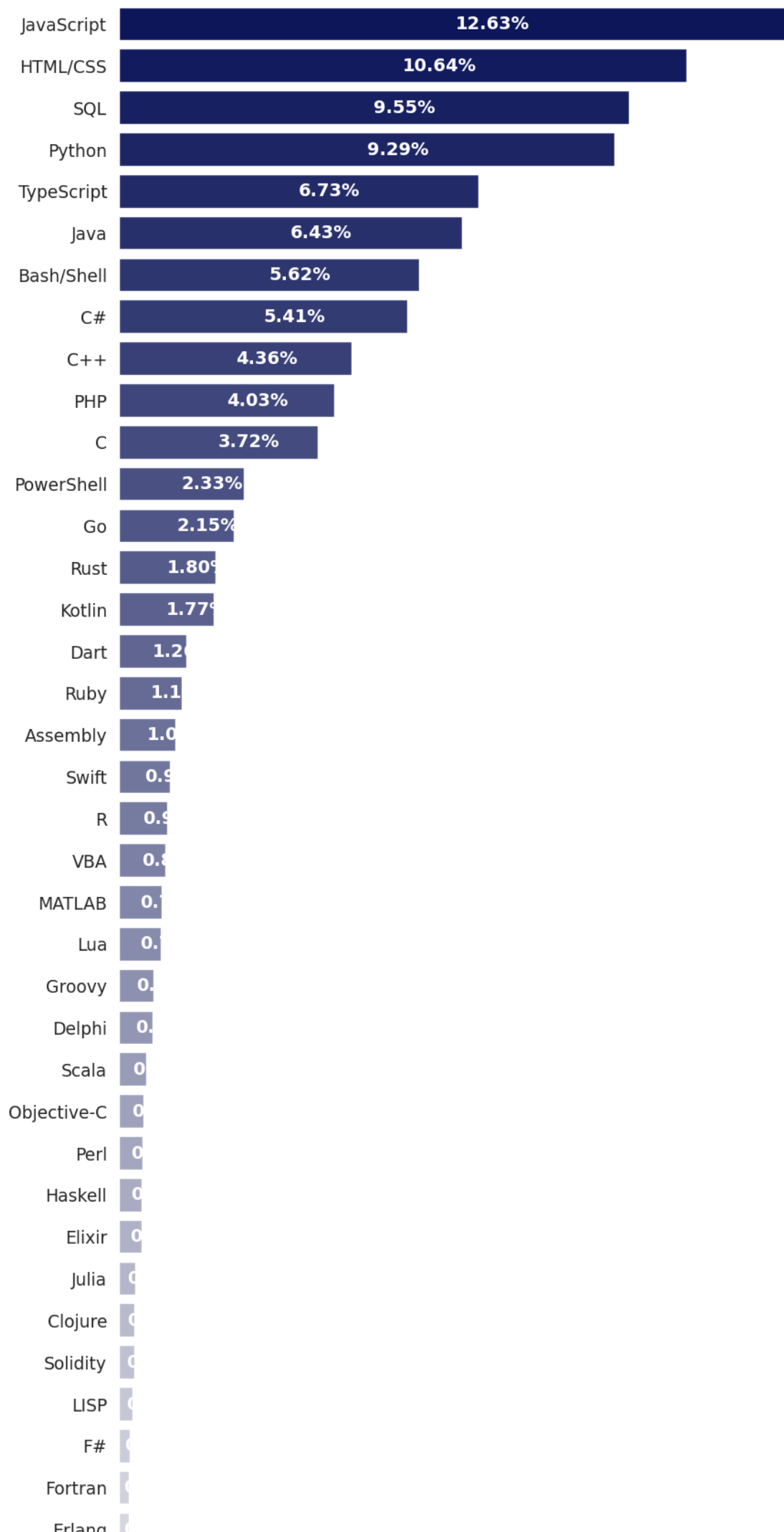
**My company has:**



Total Responses: 119166

# Which programming, scripting, and markup languages have you done extensive development work in over the past year

```
In [128]:  1  languages = colum_expand(survey_df.LanguageHaveWorkedWith).sort_values(ascen
           2
           3  s = 'Which following languages have you worked with?'
           4
           5  custom_plot(languages, plot_height=25,plot_width=10, title=s,color = 'light:
```

## Which following languages have you worked with?

| Language | Percentage |
|---|---|
| JavaScript | 12.63% |
| HTML/CSS | 10.64% |
| SQL | 9.55% |
| Python | 9.29% |
| TypeScript | 6.73% |
| Java | 6.43% |
| Bash/Shell | 5.62% |
| C# | 5.41% |
| C++ | 4.36% |
| PHP | 4.03% |
| C | 3.72% |
| PowerShell | 2.33% |
| Go | 2.15% |
| Rust | 1.80% |
| Kotlin | 1.77% |
| Dart | 1.2 |
| Ruby | 1.1 |
| Assembly | 1.0 |
| Swift | 0.9 |
| R | 0.9 |
| VBA | 0.8 |
| MATLAB | 0. |
| Lua | 0. |
| Groovy | 0. |
| Delphi | 0. |
| Scala | 0 |
| Objective-C | 0 |
| Perl | 0 |
| Haskell | 0 |
| Elixir | 0 |
| Julia | 0 |
| Clojure | 0 |
| Solidity | 0 |
| LISP | 0 |
| F# | 0 |
| Fortran | 0 |
| Erlang | |

APL
COBOL
SAS
OCaml
Crystal

**Total Responses: 367821**

In [ ]:    1

In [ ]:
```python
 1  survey_df.WebframeHaveWorkedWith      # django flask
 2  survey_df.WebframeWantToWorkWith
 3
 4  survey_df.LanguageWantToWorkWith
 5
 6  survey_df.DatabaseHaveWorkedWith
 7  survey_df.DatabaseWantToWorkWith
 8
 9  survey_df.PlatformHaveWorkedWith
10  survey_df.PlatformWantToWorkWith
11
12  survey_df.MiscTechHaveWorkedWith
13  survey_df.MiscTechWantToWorkWith
14
15  survey_df.ToolsTechHaveWorkedWith
16  survey_df.ToolsTechWantToWorkWith
17
18  survey_df.CompTotal                   # annual income
19  survey_df['OpSysPersonal use']        # operating system
```

# E N D