# *Caption Generator for Images powered by Watson Visual Recognition*

*Project report submitted by*

*Dr. B. Ramani,*

*Sri Sivasubramaniya Nadar College of Engineering*

**Introduction:**

Image Captioning refers to the process of generating textual description from an image, based on the objects and actions in the image. The ability to recognize image features and generate accurate, syntactically reasonable text descriptions is important for many tasks in computer vision. For example:

(1) Self-driving cars:

Automatic driving is one of the biggest challenges and if we can properly caption the scene around the car, it can give a boost to the self-driving system.

(2) Aid for visually challenged:

This can help us create a product for the blind which will guide them to travel on the roads without the support of anyone else. This can be implemented by first converting the scene into text and then the text to voice.

CCTV cameras are everywhere today, but along with viewing the world, if we can also generate relevant captions, then we can raise alarms as soon as there is some malicious activity going on somewhere. This could probably help reduce some crime and/or accidents.

**Proposed Solution :**

This project aims at building an application which takes input as image analyses it and generate the captions in the form of speech. To achieve this, IBM Services like node-red service is used to build a web UI where user uploads a picture. This picture is analyzed by Watson visual recognition service and the analyzed description is then converted in to speech using text to speech service.

**Workflow of the project:**

The project involves the 3 processes mentioned below.

Process 1: Create IBM Academic Initiative Account

      Step 1: IBM Academic Initiative Account

Step 2: Create Node-red Starter Service

Step 3: Create Visual recognition Service

Step 4: Create Text to speech service

Process 2: Build basic node red flow

Step 1: Install browser-utils

Step 2: Install dashboard nodes

Step 3: Nodes to be used

Process 3: Create Node-red Dashboard

Step 1: Insert template node to create Basic html Page
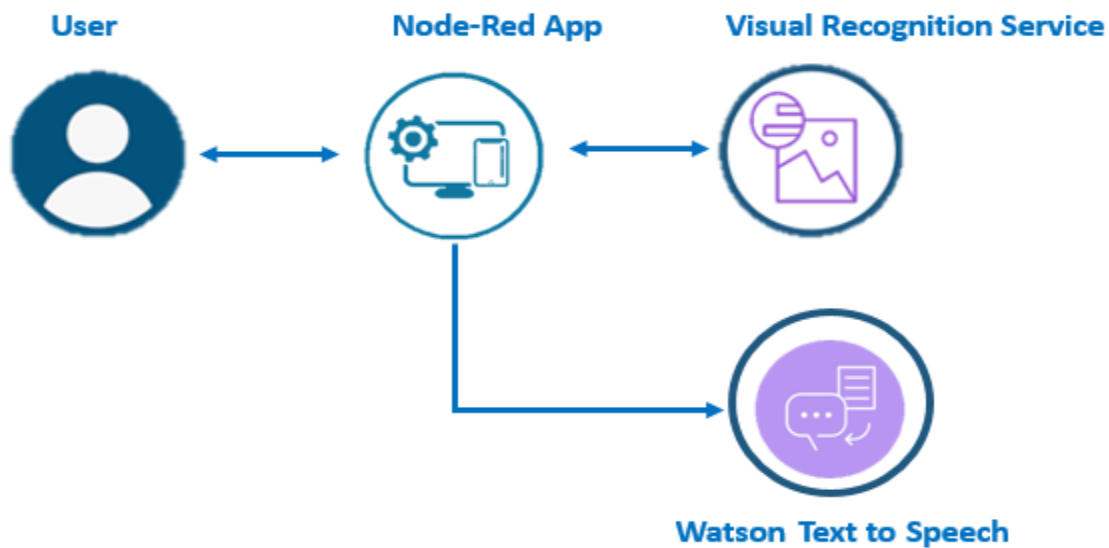
Step 2: Deploy and Open the page



Fig. 1 Work flow of the caption generator

**Implementation of the project:**

Visual Recognition Service:

The project is implemented by using the pre-built model in Visual Recognition service, which generates class keywords that describe the image. Given any input image, the classes are identified and given a confidence score which is nothing but the probability with which each class is identified.

Text to Speech Service:

The text to speech service enables us the provide speech output for any text given to it as input.

Node-Red Service:

Once the classes in the image is classified, in order have a web UI, node-red service is used. The web UI design is designed such a way that it gets an image as input from the user. This is then analyzed by Visual Recognition service and we obtain the class and the corresponding confidence scores. The caption text to be displayed is decided based on the confidence score and this is fed to the text to speech service to read that text. Fig. 2 shows the different nodes deployed for the web UI using node-red.
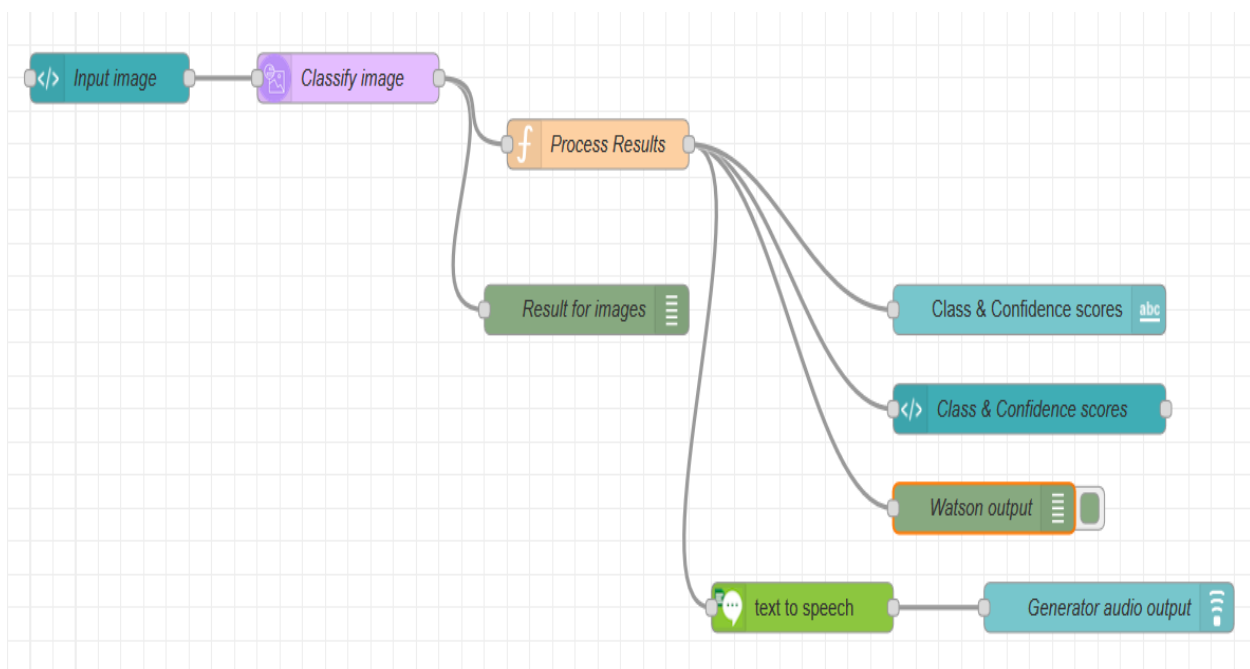


Fig.2 Nodes deployed in node-red

**Experimetal Results:**

Once the nodes are deployed, the web UI is displayed as shown in Fig. 3.
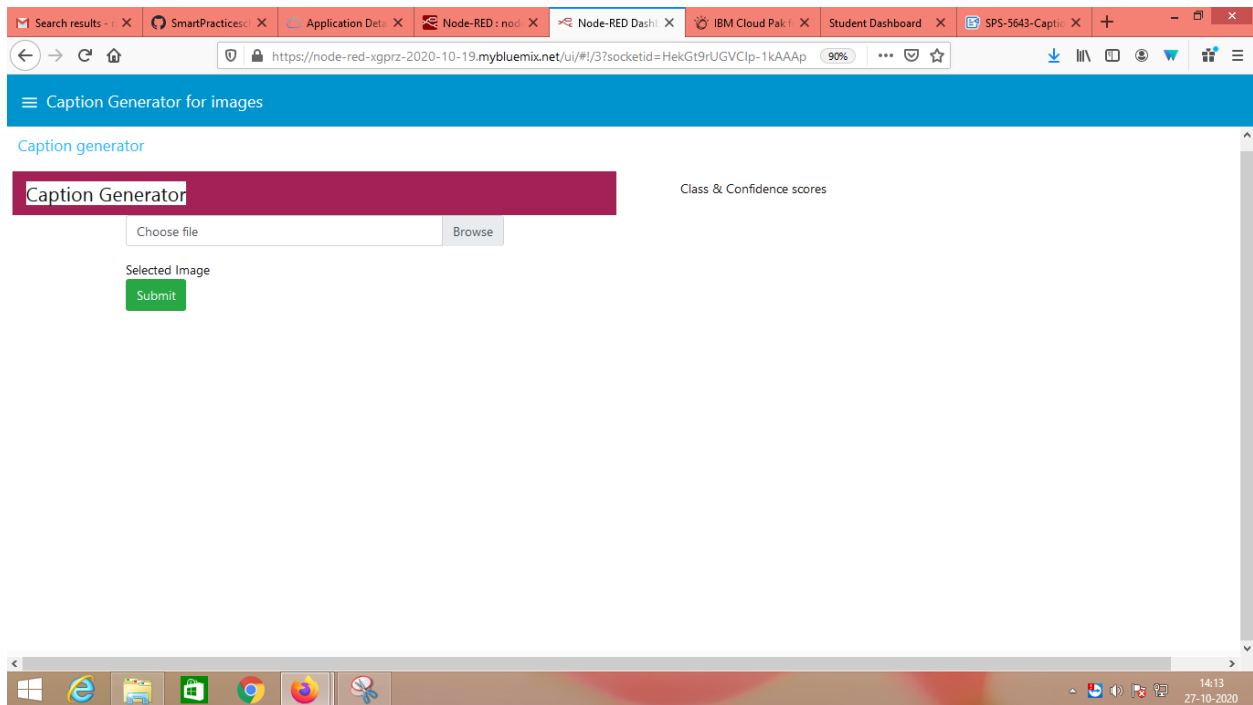
Fig. 3 Deployment of web UI using node-red service

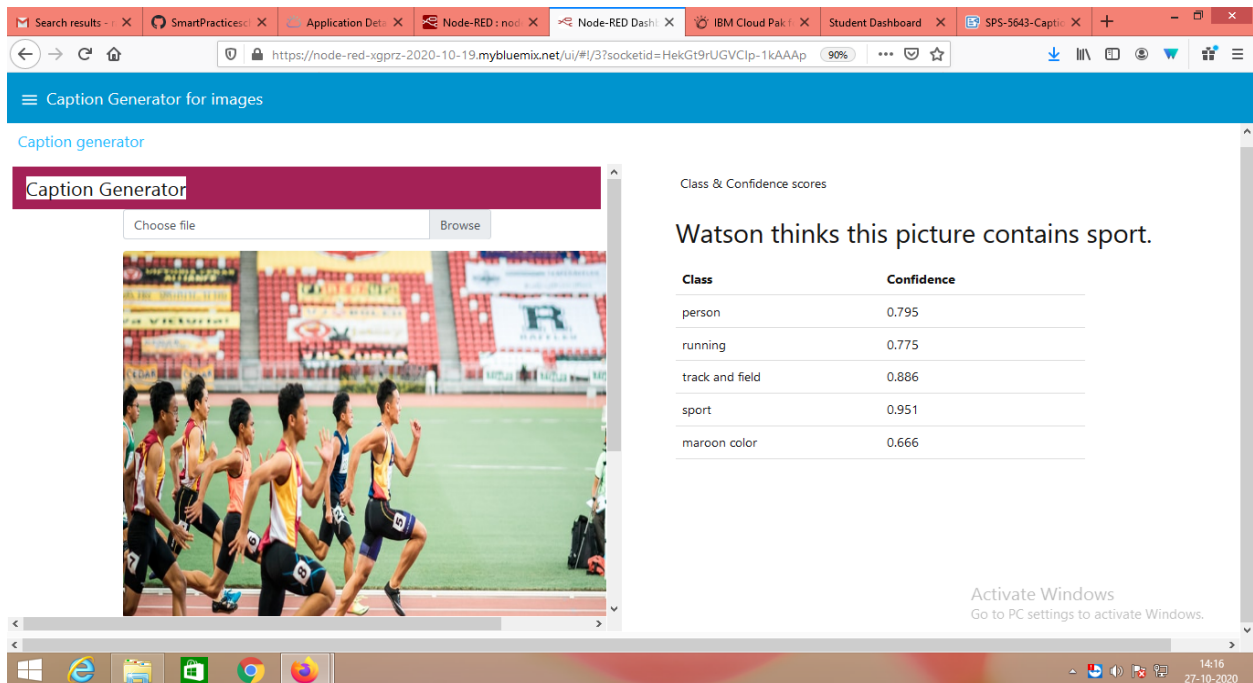By uploading different images, the results are obtained as shown in Figs. 4-6.
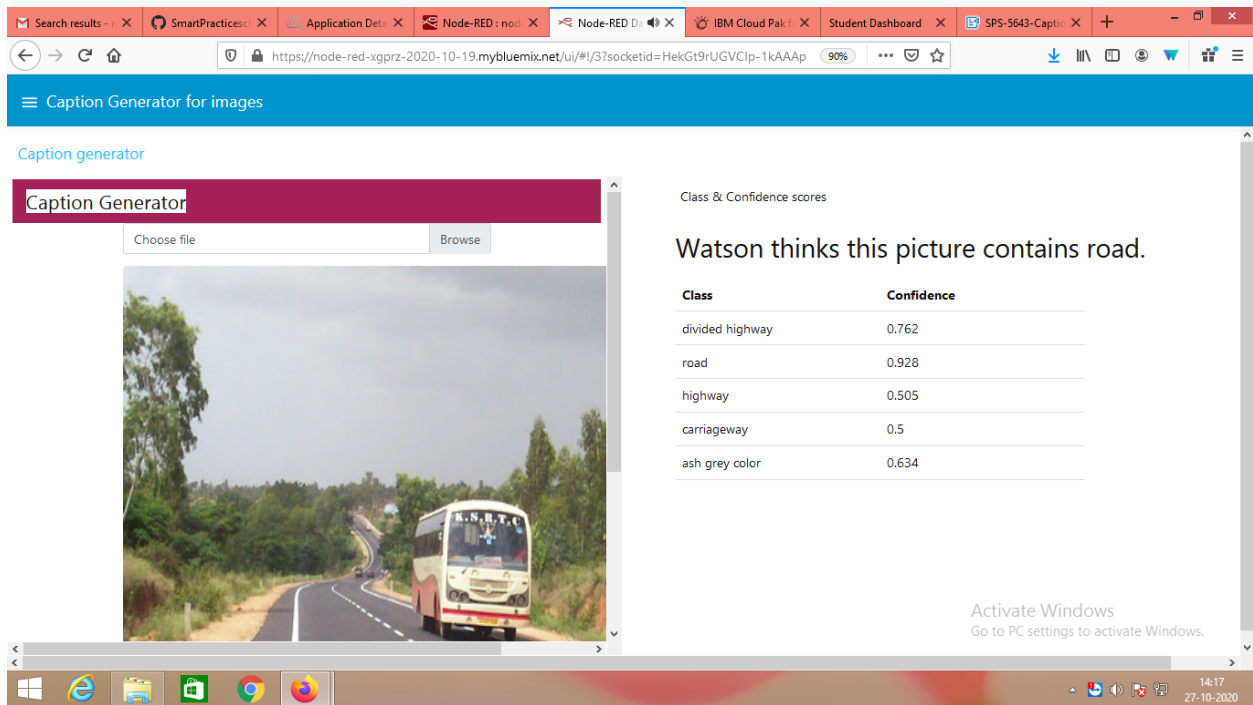


Fig. 4 Output for image 1
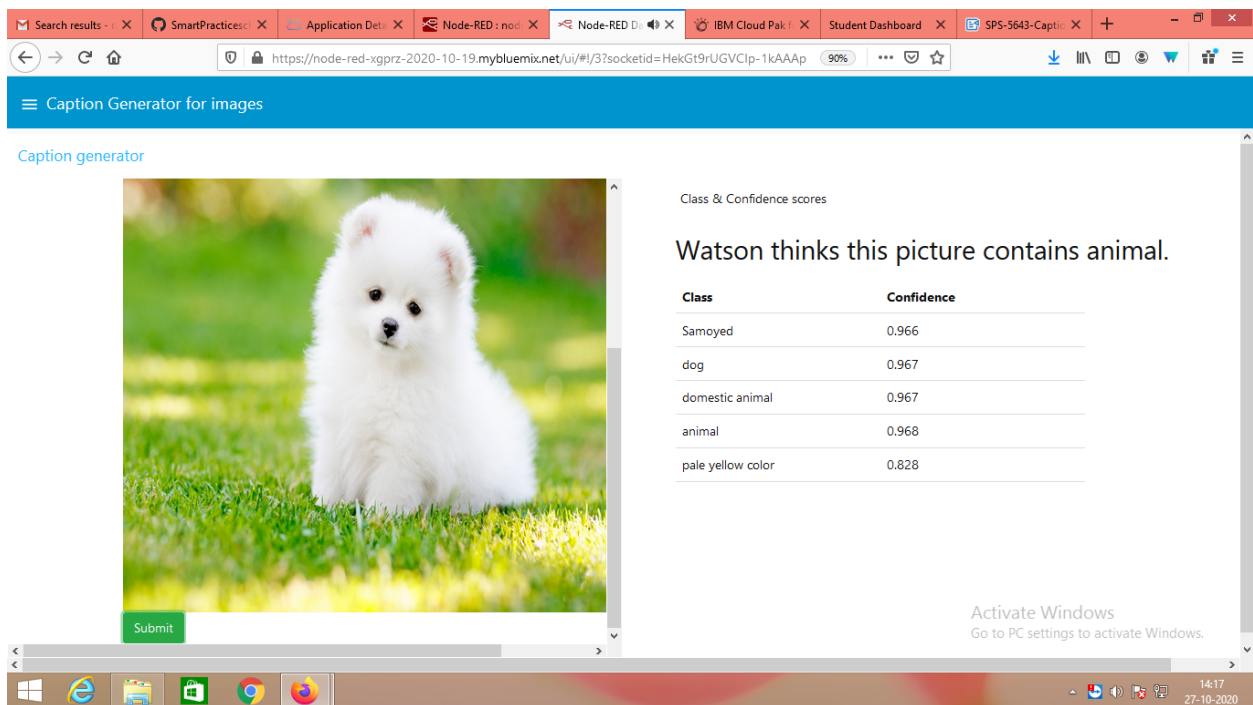
Fig. 5 Output for image 2



Fig. 6 Output for image 3

**Conclusion:**

The caption generator for images is developed using several Watson services such as visual recognition, text to speech and node-red. This caption generator can be used for several applications such as autonomous driving, aid for visually challenged, alarming the authorities in case of accidents or crime scenes.

**Future scope:**

The developed caption generator can be made effective by giving the proper caption instead of identifying a single keyword.

**Appendix**

**A.    Source code to obtain input image from user**

```
1  <!DOCTYPE html>
2  <html lang="en">
3  <head>
4    <title>Visual Recognition</title>
5    <meta charset="utf-8">
6    <meta name="viewport" content="width=device-width,
   initial-scale=1">
7    <link rel="stylesheet"
   href="https://maxcdn.bootstrapcdn.com/bootstrap/4.3.1/css/boo
   tstrap.min.css">
8    <script
   src="https://ajax.googleapis.com/ajax/libs/jquery/3.4.1/jquer
   y.min.js"></script>
9    <script
   src="https://cdnjs.cloudflare.com/ajax/libs/popper.js/1.14.7/
```

```
      umd/popper.min.js"></script>
10   <script
   src="https://maxcdn.bootstrapcdn.com/bootstrap/4.3.1/js/boots
   trap.min.js"></script>
11   <style>
12   .bg-light {
13     background-color: #A42156!important;
14           }
15           dashboard-template h4 {
16     color: #111111 !important;
17     background-color: #A42156 !important;
18     }
19
20           .custom-file {
21             margin-bottom: 14px;
22             }
23
24           .table .thead-dark th {
25             color: #fff;
26             background-color: #A42156;
27             border-color: #A42156;
28           }
29           h5{
30             text-align: center;
31               color: #EB9E30;
32             }
33             .text-center {
34               text-align: center;
35             }
```

```
36          .imgdiv
37          {
38            align:left;
39          }
40          </style>
41 </head>
42 <body>
43          <nav class="navbar navbar-expand-sm bg-light">
44            <div class="justify-content-center">
45                  <h4 class="text-center">Caption Generator</h4>
46            </div>
47          </nav>
48          <br><br>
49    <div class="container">
50            <div class="row">
51                  <div class="col-sm-2">
52                  </div>
53                  <div class="col-sm-8">
54
55                              <div class="custom-file">
56
57                                  <input type="file" name="pic"
   accept="image/*" onchange="readURL(this);" class="custom-file-input"
   id="customFile">
58
59                                  <label class="custom-file-label"
   for="customFile">Choose file</label>
60                              </div>
61                              <br>
```

```
62                              <div class="imgdiv">
63                                  <img src="#" id="blah" class="rounded"
    alt="Selected Image">
64                              </div>
65                              <!--<md-button
    ng-click="send({payload:action()})">
66                                  Predict
67                              </md-button>-->
68                              <button type="submit"
    ng-click="send({payload:action()})" class="btn btn-success">Submit</button>
69
70
71              </div>
72
73 </body>
74 </html>
75 <script>
76 var x="";
77     function readURL(input) {
78         if (input.files && input.files[0]) {
79             var reader = new FileReader();
80
81             reader.onload = function (e) {
82                 $('#blah')
83                     .attr('src', e.target.result)
84                     .width(700)
85                     .height(700);
86             };
87
```

```javascript
88                reader.readAsDataURL(input.files[0]);
89                x= input.files[0]
90
91
92          }
93      }
94      function getdata(data)
95          {
96              var html = '';
97                if(data != 0)
98                    {
99                              $.each(data, function(i){
100                                var row = data[i];
101                                console.log(row);
102                                html += '<tr>';
103                                html += '<td>';
104                                html += row.class;
105                                html += '</td>';
106                                html += '<td>';
107                                html += row.score;
108                                html += '</td>';
109                                html += '</tr>';
110                    });
111                    }
112                else
113                        html+="<div>No Data</div>";
114                $('#scoretable').html(html);
115          }
116      (function(scope) {
```

```
117          scope.$watch('msg.payload', function(data) {
118                console.log('Position 2');
119                console.dir(data);
120                getdata(data);
121
122          });
123       })(scope);
124  // or overwrite value in your callback function ...
125  this.scope.action = function() { return x; }
126
127  </script>
```

**B.    Source code to process the results obtained from Visual Recognition service**

```
1  if (typeof msg.result == 'undefined') {
2      return null;
3  }
4
5  // Text Extraction
6  if (typeof msg.result.images[0].text != 'undefined') {
7      var image_text = msg.result.images[0].text;
8      msg.payload = image_text;
9      msg.template = image_text;
10     if( image_text.length >0 ) {
11        msg.template= "Watson found the words: "+image_text;
12     }
13     return msg;
14 }
```

```
15
16 var c_id = 0;
17 var say = "";
18 var item;
19
20 for ( c_id=0; c_id <
   (msg.result.images[0].classifiers.length); c_id++ ){
21        for( i =0;
   i<(msg.result.images[0].classifiers[c_id].classes.length);
   i++ ){
22        var object =
   msg.result.images[0].classifiers[c_id].classes[i].class;
23            }
24     var bestItem = 0;
25     var itemScore = 0;
26     for( i =0;
   i<(msg.result.images[0].classifiers[c_id].classes.length);
   i++ ){
27        var object =
   msg.result.images[0].classifiers[c_id].classes[i].class;
28        if ( !object.includes("color") ) {
29          if(
   msg.result.images[0].classifiers[c_id].classes[i].score >
   itemScore){
30              bestItem = i;
31              itemScore =
   msg.result.images[0].classifiers[c_id].classes[i].score;
32          }
33        }
```

```
34      }
35      item =
   msg.result.images[0].classifiers[c_id].classes[bestItem].clas
   s;
36      say = say + " Watson thinks this picture contains " +
   item +".";
37 }
38 msg.payload =  say;
39
40 var picInfo = msg.result.images[0].classifiers[0].classes;
41 var arrayLength = picInfo.length;
42
43 msg.template="<style>h4 { text-align: center; margin: 10px;
   }";
44 msg.template=msg.template+"table {    width: 500px;
   margin-top: 10px; }";
45 msg.template=msg.template+"th, td { padding: 8px; text-align:
   left; border-bottom: 1px solid #ddd; background-color:
   #FFFFFF; width: 50%;}";
46 msg.template=msg.template+".classifier {background-color:
   rgb(85,150,030);text-align: center;}";
47 msg.template=msg.template+".title {
   background-color:LightGrey;}</style>";
48
49 msg.template=msg.template+"<h2>"+say+"</h2><table
   span=100%><tr><th>Class</th><th>Confidence</th></tr>";
50 for (var i = 0; i < arrayLength; i++) {
51   msg.template = msg.template + "<tr><td>" + picInfo[i].class
   + "</td><td>" + picInfo[i].score + "</td></tr>";
```

```
52 }
53 msg.template = msg.template + "</table>"
54
55
56 return msg;
```

**Node-red URL for web UI:**

https://node-red-xgprz-2020-10-19.mybluemix.net/ui/#!/3?socketid=ud4MRnsZFGdStlDmAAAs