# Project Proposal

*Chromilo Amin | Friday, March 26, 2021 | Bertelsmann Technology Scholarship program*

## Data Labeling Approach

| | |
|---|---|
| **Project Overview and Goal**<br><br>What is the industry problem you are trying to solve? Why use ML in solving this task? | The goal is to build a product that will help doctors quickly identify cases of pneumonia in children. We hope to train ML to correctly identify with high confidence level the presence of pneumonia from a set of images. The images will have a mix of healthy and abnormal lungs. In summary,<br>• help flag serious cases,<br>• quickly identify healthy cases,<br>• and, generally, act as a diagnostic aid for doctors. |
| **Choice of Data Labels**<br><br>What labels did you decide to add to your data? And why did you decide on these labels vs any other option? | Cloudiness_present – this is a yes-no selection<br><br>Diapragm_present – this is a yes-no selection<br><br>Heart_present - this is a yes-no selection<br><br>Cloudy_types this is a checkbox selection asking annotator to indicate whether the cloudiness is a big patch or various small patches, and on which side of the lung<br><br>Confidence_level – this is a number from 1-3 with 1 being not confidence, 2 as somewhat confident, and 3 the most confident in the answers.<br><br>Confidence_detail – this is a text field that comes up if the annotator picked 2 for confidence level above.<br><br>The pros for the above approach:<br>1- Simple yes/no questions avoids grey areas;<br>2- Unknown categories show up in confidence level, again kept simple with 1-3 only options with 2 being unsure.<br>3- Free text field gives annotators an area to voice their reason for being unsure on their answers.<br><br>The cons for the above approach:<br>1- Being overly simple especially in confidence_level ranking could result in further followups to determine if they are more or less sure about how unsure they are, i.e. expanding to 1-5 or even 1-10 choices may be improve ranking.<br>2- Lack of free text fields for each data label could hinder question improvements and increases time to surveys at the end of the job. |

# Test Questions & Quality Assurance

| | |
|---|---|
| **Number of Test Questions**<br><br>Considering the size of this dataset, how many test questions did you develop to prepare for launching a data annotation job? | I picked 3 images with abnormal readings and 2 with healthy readings for a total of 5-6 test question samples. |
| **Improving a Test Question**<br><br>Given the following test question which almost 100% of annotators missed, statistics, what steps might you take to improve or redesign this question? | <br><br>I will make sure this test ID uses the validates="required" attribute to ensure the annotators don't miss it.<br>I will make sure to provide tips why an answer is correct, incorrect, or unknown.<br>For the questions, I will include a free text comments field for annotators to use to explain and provide a reason for their answer.<br>I will review the images provided to the annotators to ensure I can classify them myself if I was presented the same questions I am asking them.<br>I will review with all the annotators the reason why they missed the question and interview an expert annotator to get help on improving the test question. |
| **Contributor Satisfaction**<br><br>Say you've run a test launch and gotten back results from your annotators; the instructions and test questions are rated below 3.5, what areas of your Instruction document would you try to improve (Examples, Test Questions, etc.) | <br><br>I think that reviewing the questions with subject matter experts may be biased itself because they are experts in the domain. I might review with non-SMEs as well.<br><br>The annotators should come from different geographies or backgrounds or cultures to normalize biases. |

# Limitations & Improvements

| | |
|---|---|
| **Data Source**<br><br>Consider the size and source of your data; what biases are built into the data and how might the data be improved? | The dataset contains 117 rows of images. To ensure non-biased results, the number of questions asked should evenly contain healthy, unhealthy, and inconclusive data sets. I used 1 test question for every 19 data points so that is about 6 questions needed.<br><br>I would suggest using more data points for images. 118 may not be representative of all children patients. For example, I would use xrays from other hospital databases. |
| **Designing for Longevity**<br><br>How might you improve your data labeling job, test questions, or product in the long-term? | I would continually survey the human annotators to query their satisfaction with the test questions, ease of completing them, areas of improvement, and providing some incentive for helping improve the data labeling job like offering gift cards.<br><br>I would submit a quick survey after the annotation jobs to get feedback. Responses as part of the test questions are different from post surveys which could garner more feedback about why annotators answered the way they did.<br>I would use expert annotator as my control data when comparing with non-experts. I do this because I think expert annotators are biased. |

# Appendix A – CML code

```
<div class="row-fluid">
  <div class="span6">
    <img src="{{hosted_image}}"/>
  </div>
  <div class="span6">
    <cml:radios name="cloudiness_present" label="Do you see any areas of abnormal cloudiness/opacity in the lung?"
validates="required">
      <cml:radio value="yes" label="Yes"/>
      <cml:radio value="no" label="No"/>
    </cml:radios>
    <cml:radios name="diapragm_present" label="Do you see a diaphragm shadow?" validates="required">
      <cml:radio value="yes" label="Yes"/>
      <cml:radio value="no" label="No"/>
    </cml:radios>
    <cml:radios name="heart_present" label="Do you see the heart?" validates="required">
      <cml:radio value="yes" label="Yes"/>
      <cml:radio value="no" label="No"/>
    </cml:radios>
    <cml:checkboxes name="cloudy_types" label="What area of the lung is opaque?" validates="required" only-
if="cloudiness_present:[yes]" exact="true">
      <cml:checkbox value="lhs" label="Big patch on left lung"/>
      <cml:checkbox value="rhs" label="Big patch on right lung"/>
      <cml:checkbox value="several" label="Several small cloudy areas on both sides"/>
    </cml:checkboxes>
    <cml:select name="confidence_level" label="How confident are you?" validates="required" only-
if="cloudiness_present:[yes]" exact="true">
      <cml:option label="1-not confident" value="1"/>
      <cml:option label="2-somewhat confident" value="2"/>
      <cml:option label="3-very confident" value="3"/>
    </cml:select>
    <cml:textarea name="confidence_detail" label="Confidence Level detail:" validates="required" only-
if="confidence_level:[2]" exact="true"/>
  </div>
</div>
```
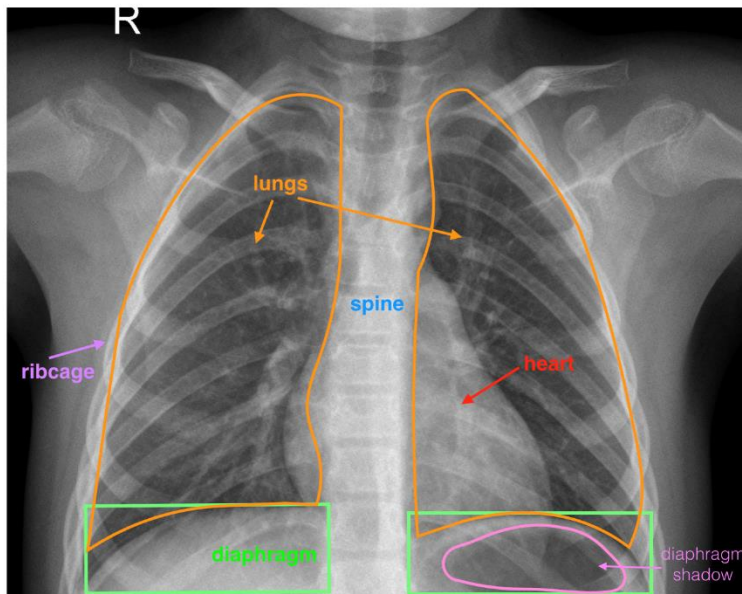
## Overview
Help us determine if the images show a healthy lung or if the patient has pneumonia.

---

## Steps
1. Examine the image and properly identify these annotated areas in the image.

* lungs
* spine
* heart
* ribcage
* diaphragm (below the lungs)



2. Check the appropriate box if the image shows clear lungs without any areas of abnormal cloudiness/opacity.

3. Check the appropriate box if the image shows the diaphragm shadow.

4. Check the appropriate box if the image shows the heart.

5. Indicate your confidence level in making your above choices from 1-3.

6. If you are somewhat confident (you picked confidence level 2 above), then type your reasons.
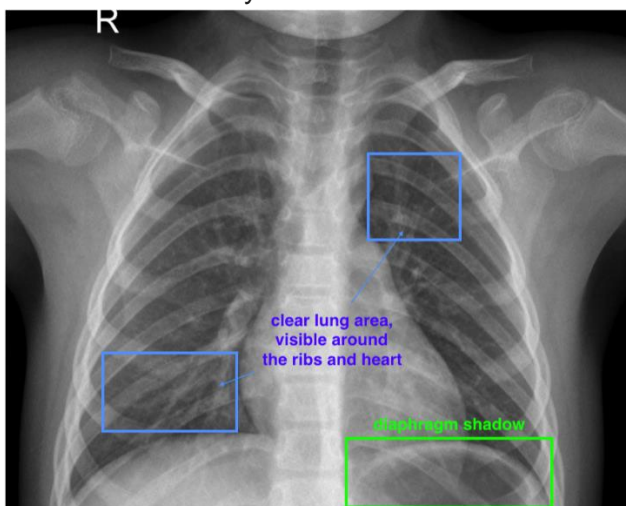
## Rules Tips
**Rules:**

---

- Categorize cloudiness as either whole patches or small regions;
- If you are not confident in your answers, type the reason.

**Tips**:

- The R in every image indicates the right hand side. Please be careful when answering the checkbox indicating which side of the lung is cloudy.

# Examples

## Normal and healthy



clear lung area, visible around the ribs and heart

diaphragm shadow

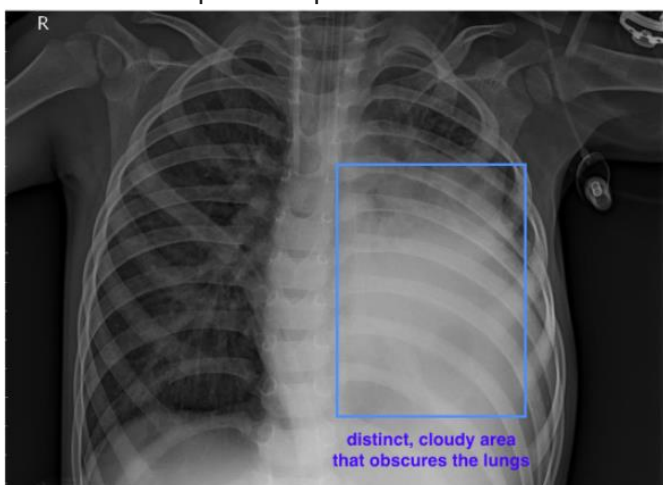**Do you see any areas of abnormal cloudiness/opacity in the lung?** (required)
○ Yes
◉ No

**Do you see a diaphragm shadow?** (required)
◉ Yes
○ No

**Do you see the heart?** (required)
◉ Yes
○ No

## Abnormal with possible pneumonia



distinct, cloudy area that obscures the lungs

**Do you see any areas of abnormal cloudiness/opacity in the lung?** (required)
◉ Yes
○ No

**Do you see a diaphragm shadow?** (required)
○ Yes
◉ No

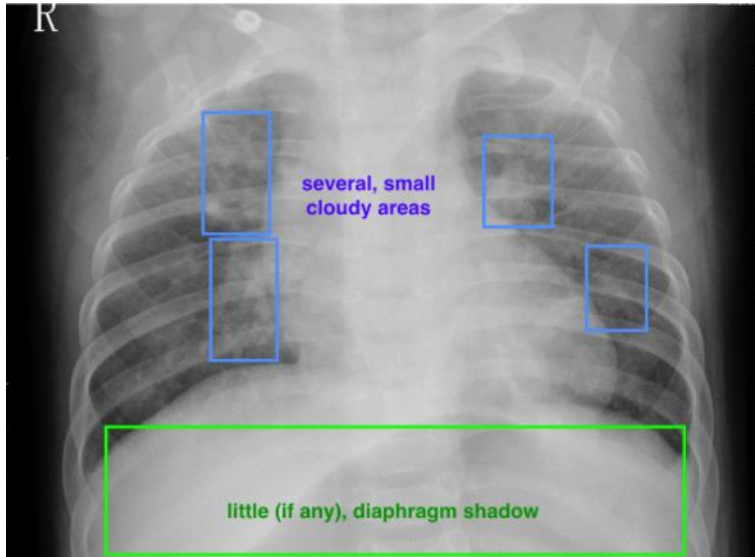**Do you see the heart?** (required)
○ Yes
◉ No

**What area of the lung is opaque?** (required)
☑ Big patch on left lung
☐ Big patch on right lung
☐ Several small cloudy areas on both sides

**How confident are you?** (required)
[ 3-very confident ⌄ ]

## Abnormal with possible pneumonia



**Do you see any areas of abnormal cloudiness/opacity in the lung?** (required)
- ◉ Yes
- ○ No

**Do you see a diaphragm shadow?** (required)
- ○ Yes
- ◉ No

**Do you see the heart?** (required)
- ◉ Yes
- ○ No

**What area of the lung is opaque?** (required)
- ☐ Big patch on left lung
- ☐ Big patch on right lung
- ☑ Several small cloudy areas on both sides

**How confident are you?** (required)

| 2-somewhat confident ⌄ |
|---|

**Confidence Level detail:** (required)

I see various small patches on both lungs.