**ChatGPT**

# Predictive Quality Monitoring System for Continuous Steel Casting (Technical Design)

## Introduction and Context

In modern steel manufacturing, **continuous casting** is a critical process where molten steel is solidified into slabs or billets for further processing [1]. Ensuring high product quality from the caster is challenging due to the complex, chaotic nature of the process with **hundreds of variables affecting outcomes** [2]. KISS Technologies (Holland, MI) specializes in data-driven solutions for continuous casting, aiming to "take the guesswork out of casting" by using sensor data and analytics to improve quality [2] [3]. The goal of this project is to design a **machine learning (ML) system that predicts the quality of steel output** in real-time based on metrics collected from casting machines. This technical specification outlines a detailed design for a proof-of-concept (PoC) demo system, including data sources, model choices, system architecture, and implementation steps. The focus is on **predicting quality defects** (e.g. surface cracks, shape deformities, inclusions) from process sensor data (e.g. temperatures, speeds, flows), with an initial prototype using synthetic data and a roadmap toward real production integration.

**Scope and Objectives:** The system will ingest streaming sensor data from the continuous caster, perform real-time ML inference to predict quality outcomes (such as a "defective" vs "good" classification or a quality score), and alert or log predictions for operators. The demo will be built with **Python** and **PyTorch**, and designed for rapid development (short timeline) while following best practices for production-quality code. Key objectives include:
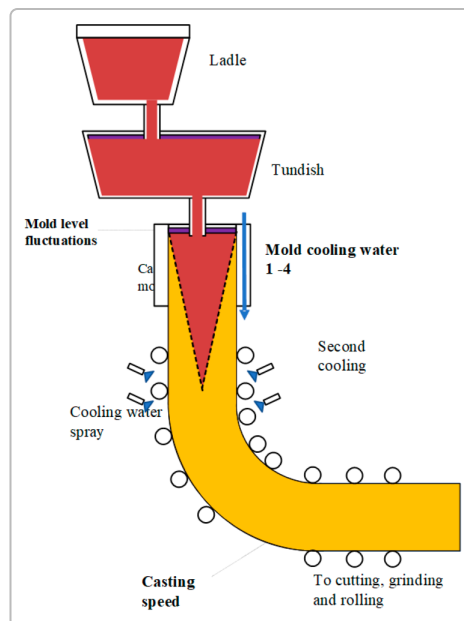
- Demonstrating domain understanding of continuous casting and relevant quality metrics.
- Evaluating different modeling approaches (from classical ML to deep learning) for predicting steel quality.
- Implementing a data pipeline that can be extended to real plant data (simulate with synthetic data initially).
- Ensuring the design considers real-time deployment, scalability, and maintainability for a production environment.

## Domain Background: Continuous Casting Process & Quality Factors

Continuous casting (strand casting) converts molten steel from a **ladle** (large vessel) through a **tundish** (distributor) into a water-cooled **mold**, where a solid shell forms, and the steel strand is continuously withdrawn, bent, and cooled in a curved path [1]. The strand is then cut into slabs or billets for rolling and finishing. The quality of the cast steel can be affected by many process factors: - **Thermal conditions:** Steel temperature (superheat above melting point) and cooling rates (mold cooling water, spray cooling) influence solidification. Improper temperature control can lead to internal cracks or segregation if too low, or surface defects if too high [4] [5]. - **Mechanical/operational parameters:** Casting speed, mold oscillation frequency and amplitude, mold level stability, and electromagnetic stirring all affect how the steel solidifies. For example, excessive casting speed or mold level fluctuations can entrain slag inclusions

or cause shape deformations [6] [7] . - **Material properties:** Steel chemical composition (carbon content, etc.), alloying additions, and molten steel cleanliness (inclusion content) alter viscosity and solidification behavior [8] [9] . Variations in composition can change the final microstructure and quality of the slab, requiring process adjustments. - **Environmental factors:** Ambient temperature, humidity, and air flow in the caster environment can affect cooling efficiency [10] . For instance, high ambient temperatures might reduce cooling, increasing the risk of certain surface defects [11] . - **Equipment condition:** The state of machinery (oscillators, rollers, cooling pumps, etc.) can indirectly impact quality. Worn oscillation mechanics or clogged cooling nozzles might create non-uniform conditions leading to defects [12] . Vibration or temperature sensors on equipment can indicate issues that correlate with quality problems (or signal need for maintenance).

**Quality Defects of Interest:** Common continuous casting defects include **surface cracks**, **internal cracks**, **corner cracks, shape deformities** (e.g. rhomboidity of billets where cross-section becomes rhombus-shaped [4] ), and **non-metallic inclusions** (like slag entrapments). These defects can cause downstream problems in rolling and reduce product quality. For example, rhomboidity can arise from uneven cooling or mold issues and is measured by the difference in slab diagonal lengths [4] . Slag inclusions often result from mold level fluctuations or improper flow, leading to embedded particles in the steel [13] [14] . A predictive system aims to **forecast the occurrence or severity of such defects** from sensor data, allowing operators to correct process parameters in real-time (e.g. adjust cooling or speed) to prevent defects [15] .



*Schematic of a continuous casting process for steel, from ladle (top) to solidified strand (bottom). The mold creates a solid shell as steel flows through it, with mold level control at the top. Mold cooling water (labeled 1–4 on sides) and secondary cooling spray regulate solidification. Casting speed controls strand withdrawal. These process parameters are continuously monitored by sensors and profoundly affect final steel quality.* [16] [6]

# Data Sources and Collection Strategy

A successful predictive model requires integrating diverse data from the continuous casting operation. Below we identify key data sources and how they would be collected in a modern steel plant setup (and simulated for the PoC):

- **Process Parameters:** Real-time sensor readings that characterize the casting operation. This includes **casting speed (strand withdrawal rate)**, molten steel **temperatures** (e.g. in tundish and mold), **coolant water temperatures** and **flow rates**, **mold oscillation** settings (frequency, stroke), **mold level** (steel height in mold) and its fluctuations, and mold **friction or stopper rod position** (indicating flow control). These are typically measured via thermocouples, flow meters, level sensors, and oscillation encoders on the caster [16] . In the PoC, we will simulate these as time-series data streams (e.g., generating numeric values at 1 Hz for each parameter, mimicking real sensor feeds).

- **Material Properties:** Data describing the molten steel itself, such as **chemical composition** of the heat (from lab analysis or upstream sensor if available), **steel grade**, **inclusions count/cleanliness**, and **superheat** (temperature above liquidus) [8] . For demo purposes, we can assign synthetic composition IDs or categories and nominal superheat values. These properties may be static per casting sequence (each heat of steel) rather than changing rapidly, so they can be provided as context data at the start of a cast simulation. Monitoring these is crucial because composition variations can alter optimal casting parameters and quality outcomes [9] .

- **Environmental Data:** Ambient conditions around the caster, especially **air temperature**, **humidity**, and possibly **cooling water inlet temperature/pressure** if affected by environment [10] . We can simulate these as slowly varying parameters. In reality, this data might come from plant environmental sensors or weather data integrated with the control system. These factors help explain certain quality variations (e.g. hotter days might require adjusting cooling water flow to maintain the same cooling effect [11] ).

- **Equipment Sensor Data:** Machine condition monitoring data, such as **vibration readings** on oscillation drive or support rollers, **hydraulic pressure** in the mold oscillation system, and **motor currents/temperatures** [12] . These indicate if equipment is operating normally. In a data-driven system, such metrics can be used for **predictive maintenance** and also correlated with quality (e.g., abnormal vibrations might cause oscillation deviations affecting the strand surface). For the PoC, we might include a simplified indicator (e.g., an "oscillator health" value or a binary "maintenance needed" flag) to demonstrate how equipment status could be factored in.

- **Production and Quality Outputs:** Finally, the **label or target** for our ML model comes from production data: measurements of the final product quality. This could be **binary labels** (defective vs. acceptable slab) or **numeric quality metrics** (e.g., crack count, size of deformity, or an overall quality score). In real plants, quality is assessed by inspections (surface scanners, ultrasonic tests for internal defects, etc., or downstream rejection/acceptance info) [17] [7] . For our synthetic dataset, we will generate a quality outcome for each cast sequence based on the simulated process conditions (e.g., flag as "defect" if certain sensor thresholds are crossed during the cast). **Production yield data** (tons cast, downtime, etc.) can also be logged, though primarily for process improvement analytics rather than direct input to the quality model [18] .

**Data Collection Approach:** In a production setting, all the above data streams would be collected via an Industrial IoT or automation system. Sensors feed into a **Level 2 database or historian** in the steel plant. For a modern setup, one could use an **MQTT broker or Apache Kafka** to stream sensor data to analytics systems, or a time-series database (OSIsoft PI, etc.) for storage. For the **demo PoC**, we will simplify: sensor readings will be generated by a Python script (or notebooks) to simulate time-stamped data. This synthetic data generator can output CSV files (for batch training) and also emit data in real-time to a message queue or socket to emulate streaming. This ensures our design covers both **batch historical data** (for model training) and **streaming live data** (for inference). We will structure the data as follows:

- **Training Dataset:** A collection of casting runs (e.g. 1000 simulated casts) with associated time-series data and a known quality outcome label for each. Each "cast" in data could be represented by summary statistics (min, max, mean, etc. of each sensor) or by the full time-series. Initially, we may compute features from the time-series for each cast for simpler model training (to avoid extremely long sequences), but we will keep the raw sequences as well for potential deep learning models.
- **Real-Time Feed:** A continuous stream of sensor readings for a current cast. The system will accumulate these readings in a short time window (or use sequence-to-sequence modeling) to predict the quality before the cast ends. For demo, we can loop through a test dataset and feed data incrementally to the model, simulating an online prediction scenario.

Throughout the PoC, careful **data management** is crucial. We will use **Pandas** for handling batch data (CSV files), and possibly a lightweight message queue (or even a socket or shared memory) for streaming. Data normalization or scaling will be applied (e.g., z-scores or min-max scaling for each sensor) based on training-set statistics, to ensure the ML model sees normalized inputs. Additionally, we will design for **data logging**: capturing predictions and key sensor values to output (for demonstrating the results and enabling later analysis).

## System Architecture Overview

To meet the requirements of real-time prediction and future scalability, the system is designed in a **modular architecture** with distinct components for data ingestion, processing, model inference, and user interface/alerting. Below is an overview of the architecture, followed by detailed technical options for each component:

- **Data Ingestion Layer:** Responsible for collecting sensor data (or reading from the synthetic data generator). This could be implemented as a **streaming client** (e.g., a Python script that subscribes to a Kafka topic or reads from a socket) that buffers incoming sensor values. For the PoC, a simple loop reading from a simulated data source (CSV or generator function) will suffice. In production, this layer would interface with plant control systems or IoT gateways to get live data.

- **Data Processing & Feature Engineering:** This layer processes raw sensor readings into the input format needed by the ML model. Options at this stage include: (a) real-time feature computation (such as calculating rolling statistics over the last N seconds of data, or aggregating the entire sequence of a cast), or (b) feeding the raw sequence into a sequence-based model. A **sliding window** approach might be used for long casts, updating prediction as data comes in. For demo simplicity, we might process the entire sequence of a simulated cast at once to predict the final quality. We will implement data cleansing here too: handling missing values (if any sensor drops

out), smoothing noisy signals if needed, and aligning timestamps. The output of this layer is a feature vector or time-series tensor ready for the model.

- **Model Inference Engine:** The core ML model that given the processed data outputs a quality prediction. This will be built in **PyTorch**, allowing flexibility for advanced networks. We will design the model component to be **replaceable**, meaning we can plug in different model variants (for experimentation) behind a common interface. The model engine can run in two modes:

  - **Offline Training Mode:** where it reads historical data, trains a model (e.g., PyTorch model training loop or using PyTorch Lightning for structure), and validates performance.
  - **Online Inference Mode:** where a trained model is loaded and applied to incoming data to produce predictions on the fly.

We plan to implement an initial model (e.g., an LSTM-based classifier for defect prediction) and possibly a simpler baseline (like a random forest using summary features) for comparison. The inference engine will output a prediction along with a confidence or score. If classification is used, this might be a probability of "defect" vs "no defect"; if regression, it could be a predicted quality metric (like expected crack length or a quality index).

- **Alerting & Visualization:** In a real system, the predictions would be delivered to operators or to other control systems. This might involve a dashboard, an HMI (Human-Machine Interface) on the caster, or alerts when quality is predicted to be poor. For the demo, we will implement a simple console log or web dashboard that updates with the latest prediction and possibly a chart of sensor readings. For example, a small web app (using a Python web framework or a Jupyter notebook widget) could show real-time graphs of key sensors (temperature, speed, etc.) and a live indicator of predicted quality. This is optional for the PoC but can significantly "impress" by visualizing the ML system in action.

- **Data Storage & Logging:** All data and predictions will be logged for analysis. During training, we'll keep a dataset of all simulated runs. During real-time operation, we will log sensor streams and model outputs to a file or database. This is important for debugging and for demonstrating how the model could be monitored in production. In an actual plant, a time-series database or cloud storage might store this information for later audit.

**Architectural Considerations:** The system is designed to be as real-time as needed (sensor data likely at 1 Hz or faster; our approach can handle this easily since the model inference (even a deep network) will be on the order of milliseconds). If higher-frequency data (e.g., dozens of Hz) or large numbers of sensors are used, careful consideration of throughput will be needed – possibly using streaming frameworks or batch processing micro-batches of data. The demo system's modular design ensures that each piece (ingestion, processing, inference, output) can be scaled or replaced without affecting others. For instance, switching from synthetic data to a real MQTT feed would only require changes in the ingestion module, while the rest of the pipeline remains the same. This modular approach and separation of concerns reflect production-quality software design.

# Modeling Approach: Options Analysis and Recommendations

**Problem Formulation:** We frame the quality prediction as a **supervised learning problem**. Each casting sequence (or a defined time window of casting) is an instance with input features (sensor time-series) and an output label or value (quality outcome). We consider two formulations: - A **classification** approach (predict if the cast will be "defective" or "acceptable" in terms of quality), which is suitable if the primary goal is to catch bad casts. - A **regression** approach (predict a numerical quality score or defect severity metric). This could provide more nuanced prediction (how bad the quality might be), but requires a well-defined quality metric.

For initial implementation, we recommend a **binary classification** (defect vs no defect), as it is simpler and aligns with common industry practice to detect out-of-spec product. This can later be extended to multi-class (for different defect types) or regression if needed.

**Model Type Options:** We evaluate a few candidate modeling techniques: 1. **Gradient Boosted Decision Trees (GBDT)** – e.g. using XGBoost or LightGBM. These models handle tabular data well and can give feature importance. They require us to first reduce our time-series to a fixed set of features per cast (such as average temperature, max temperature, standard deviation of mold level, etc.). GBDT models are fast to train and often effective on smaller datasets. *Pros:* interpretable features, easy to implement, robust to smaller data sizes. *Cons:* manual feature engineering needed, may not capture temporal patterns fully. 2. **Traditional ML (SVM, k-NN, etc.)** – Also requires feature engineering; less likely to outperform GBDT or deep learning in this context, but could be used as simple baselines. For instance, a k-NN classifier could be tried, though handling high-dimensional time-series data may be an issue. These are generally not our top recommendation given richer methods available. 3. **Deep Neural Network – Time Series Approach:** Here we use the raw or minimally processed time-series as input to a neural network: - **Recurrent Neural Networks (RNNs)** like **LSTM** or **GRU** networks are well-suited to sequence data. An LSTM-based classifier can take sequences of sensor readings (potentially multivariate sequences concatenating temperature, speed, etc. per time step) and learn temporal patterns that lead to defects. This approach was successfully used in research for continuous casting; e.g., a multi-scale CNN-LSTM model was used to detect anomalies in casting sensor data and predict inclusions [19]. *Pros:* captures time dependencies, no need to manually pick features, can potentially catch subtle temporal patterns (like a brief temperature dip or oscillation spike that precedes a defect). *Cons:* needs more data to train effectively, harder to interpret, and requires more computational effort. - **1D Convolutional Neural Networks (CNNs)** for time-series: By treating the multivariate time series like signals, a CNN can extract local temporal patterns. CNNs can also be combined with RNNs (as in the above research) to handle multi-scale patterns [19]. For example, a CNN could first process each sensor signal to extract features (like patterns of oscillation), then an LSTM layer aggregates sequence information. - **Transformer-based models**: Modern approaches like the Transformer architecture (or time-series specific adaptations) could also be considered for sequence classification. These might be overkill for our PoC, but are an option for long sequences due to their ability to capture long-range dependencies. 4. **Hybrid/Ensemble Approaches:** A combination of the above – e.g., use an ML model on summary features plus an alert from an anomaly detection model. One idea is to employ **unsupervised anomaly detection** (like an autoencoder or one-class SVM) on the sensor data to detect unusual patterns, and use that as an input feature to a main classifier. Another idea is to ensemble a tree-based model with a neural network to see if that improves robustness.

**Recommendation:** Start with a **two-stage approach**: - Implement a **baseline GBDT model** on engineered features to set an initial performance benchmark (this is quick to train and interpret). We will create features

such as: max/min/avg of temperatures, standard deviation of mold level, frequency of mold level oscillations beyond a threshold, etc., for each cast. This model can be built with scikit-learn or XGBoost. It provides feature importance which can highlight key drivers of defects (e.g., it might show mold level variability is a top predictor). - Then develop a **PyTorch LSTM-based deep learning model** that ingests the full sequence of key sensor readings. The architecture might be, for example: an LSTM layer (or stacked LSTMs) processing a sequence of shape [time_steps × features] followed by a couple of fully connected layers to output a probability of defect. We will use PyTorch to build and train this network on the synthetic dataset. We expect this model to potentially capture more complex interactions (like the timing of a temperature drop relative to a speed increase). If the dataset is small, we might need to be careful to avoid overfitting (techniques like dropout, early stopping, or even data augmentation on the time series could help).

We will compare the results of these approaches. If the LSTM outperforms the baseline, that's a good indicator of temporal patterns being important. If not, it might mean our synthetic data relationships are simple or the feature-based approach is sufficient – either way, that is a useful insight to discuss in the interview.

**Handling Imbalanced Data:** In many real cases, the occurrence of defects is relatively rare (e.g., only a few percent of casts have serious defects) [7] . This class imbalance can hurt model training (the model may always predict "good" and achieve high accuracy but be useless). Our synthetic data can be generated to either balance classes or reflect imbalance. We will employ techniques as needed: for tree models, use balanced objective or class weights; for neural nets, use weighted loss or oversample the defect cases. Additionally, **evaluation metrics** beyond raw accuracy will be used – e.g., the **ROC-AUC, F1-score, recall for defect class** – to ensure we capture the model's ability to detect the problematic cases [20] [21] .

**Feature Engineering (for non-deep models):** It's worth detailing some features we'd compute for the baseline model: - **Statistical features:** mean, median, max, min, variance of each sensor over the cast duration. - **Dynamic features:** e.g., frequency of abrupt changes or spikes (like count of temperature drops > X degrees), oscillation frequencies (from mold oscillation sensor or indirectly from mold level oscillations). - **Duration at extremes:** e.g., time spent above a critical temperature, or time mold level was unstable beyond a threshold. - **Cross-features:** combinations like casting speed * superheat (to encode cooling conditions), etc.

These will be derived from the time-series in a preprocessing script. For deep models, we feed the raw (or lightly preprocessed, e.g., normalized) time-series, possibly truncating or padding sequences to a fixed length if needed.

## Implementation Plan (Step-by-Step)

This section outlines the concrete steps to build the demo system, including both the PoC development with synthetic data and considerations for future real-data integration.

**1. Environment Setup:** Prepare a Python environment with required libraries. We will use `pandas` and `numpy` for data manipulation, `PyTorch` for model development, and possibly `scikit-learn` or `xgboost` for baseline models. Ensure that the environment can handle real-time operations – this is mostly about having the ability to run an async loop or multi-threading if needed for streaming. Also set up

any simple message passing mechanism for the stream (for example, Python's `queue.Queue` or a socket server). If time permits and it adds value, set up a lightweight **dashboard** environment (e.g., using Plotly Dash or a Jupyter notebook with live plotting) to visualize results.

**2. Synthetic Data Generation:** Develop a module to simulate continuous casting data. This involves: - Defining the **time span** of one cast (for example, one casting sequence might last say 10 minutes for a billet or longer for a sequence of slabs; we can simulate a shorter period for demo). - Generating time-series arrays for each sensor: - *Steel temperature:* e.g., start near 1550°C and gradually cool in mold, but maintain above liquidus; we can simulate a slight decline or fluctuations. - *Mold level:* simulate as roughly constant around a setpoint, with random noise and occasional larger perturbations. - *Casting speed:* perhaps constant or a controlled variation (ramping up or down). - *Cooling water flow rates:* could be constant or adjusted during cast, with random noise. - *Any equipment signals:* e.g., simulate a "vibration" signal that is low (normal) vs spikes if we introduce a fault scenario. - We then decide rules to determine quality outcome: for instance, if mold level fluctuated beyond a threshold for more than N seconds or if temperature dropped below a limit, mark that cast as "defective". These rules can be informed by domain knowledge (e.g., *"significant difference between mold diagonals leads to rhomboidity"* might correlate with uneven cooling – we can approximate by checking differences in cooling flow rates [4] [5] ). - Generate a dataset of many casts (the more the better for training; perhaps a few thousand). We will ensure there is variety (different steel grades, different random seeds for noise) so the model can generalize.

We will save this synthetic dataset to file (CSV or JSON format), with a structure like: each cast has an ID, a label, and a time-series for each sensor (possibly stored as separate time-step rows or a serialized list). We might also generate precomputed features here for the baseline model.

**3. Exploratory Data Analysis (EDA):** Before modeling, perform a quick analysis on the synthetic data to verify it behaves realistically. Plot examples of sensor readings for a "good" cast vs a "defective" cast to see if the differences are visible (for instance, maybe the defective one shows a big temperature dip). This helps sanity-check the data generation and provides visuals that could be shared in the interview. It also guides feature engineering (confirming that, say, a certain metric truly correlates with defects in our synthetic setup).

**4. Model Training - Baseline:** Implement the training of a baseline ML model: - Compute feature vectors for each cast from the raw data (if not precomputed). Use the features described earlier (ensuring to use only data that would be available by the time of prediction; e.g., since we predict at end of cast, features can span the whole cast). - Split the dataset into training and test sets (e.g., 80/20 split). Ensure class distribution is preserved or use stratified splitting if classification. - Train a model such as **XGBoost** classifier or a RandomForest on the training set. Tune hyperparameters briefly (though with synthetic data, a simple default might suffice). Use appropriate evaluation metrics (calculate accuracy, precision/recall, F1, AUC on the test set). - Examine **feature importance** from the model (many tree models provide this). We will document which parameters were most influential – this not only shows we used the model but also provides domain insight (e.g., if "casting speed variance" is top feature, that suggests controlling speed tightly is crucial). - Save the trained model (for example, pickle the model object). This baseline model can be used in parallel with the deep model during the real-time demo to show two opinions, or just as a benchmark in the report.

**5. Model Training - Deep Learning (LSTM network):** Implement a PyTorch model for sequence classification: - Prepare the data loader: We will likely pad or truncate sequences to a fixed length or use

sequence packing in PyTorch, depending on how variable our sequence lengths are. Another approach is to use the entire sequence with an LSTM and just rely on end-of-sequence output, which is fine if lengths are similar. - Define the model architecture, e.g.:

```
class QualityLSTM(nn.Module):
    def __init__(self, input_dim, hidden_dim, num_layers, num_classes):
        # define LSTM and Linear layers
```

Possibly we set `input_dim = number of sensor channels` (e.g., 5 or 6), `hidden_dim` maybe 50, `num_layers=1 or 2`, `num_classes=2` (for defect/no defect). We'll include dropout layers for regularization. - Train the network on the training set. Use a loss like **binary cross-entropy** for classification. Use an optimizer like Adam. We will train for a number of epochs (monitoring validation loss). Because our dataset is synthetic and not huge, this should train fairly quickly (we can also limit sequence length or downsample time steps if performance is an issue). - Evaluate on test set, compute the same metrics (accuracy, recall, etc.). Compare to baseline. If the deep model performs better (e.g., higher recall of defects), that's a positive. If not, it might mean the simple features captured most of the signal – but in either case we'll have a good discussion point. - Save the trained PyTorch model (state_dict or scripted module). We could also **export the model to ONNX** format if we wanted to demonstrate how it could be deployed to other environments (edge devices, etc.), but that's optional.

**6. Real-Time Inference Demo:** With models in hand, set up the pipeline to simulate live operation: - Create a loop that goes through a new simulated cast (or reuse some test cases) and feeds data incrementally. For example, at each time step (e.g., each second of simulated time), send the current sensor readings to the system. - The system's inference component should accumulate these readings. There are two strategies: a) **Event-driven prediction:** Continuously update the prediction as each new data point arrives (for instance, we could maintain a running prediction using the LSTM's state, or simply re-run the model on the whole sequence seen so far at intervals). b) **End-of-cast prediction:** Only issue a prediction after the cast is complete (which in real life means the prediction comes just in time to inform the next cast or to decide if this cast's output needs scrapping/rework). Ideally we want to show *early warning*, so approach (a) is more impressive – e.g., "halfway through the cast, the model predicts with 90% confidence that a defect is forming". - Implement a simple version of (a) for demo: perhaps every fixed interval (say every 30 seconds of simulated time), compute features or run the model on data up to that point, and output a predicted probability. This can be visualized as a time series of "defect risk" that hopefully rises if a defect is indeed going to occur. - Ensure that latency is low. Both our baseline and LSTM model should be extremely fast for one instance (milliseconds), so the main delay is reading/waiting for data. For demo simplicity, the loop can just have a small sleep to mimic real-time. - Connect the output to a user interface element. If using a console, just print the predictions. If using a GUI, update a plot or indicator. For instance, we could show a live chart: sensors on one subplot, predicted defect probability on another. - This step showcases the **working skeleton for real-time integration**: It mimics how in a plant, the model would run continuously and provide ongoing quality assessments.

**7. Evaluation and Iteration:** After building the demo, we'll prepare an analysis of results: - Summarize the performance of the model on synthetic data (e.g., "the model achieved 95% accuracy in classifying casts, detecting 90% of defective cases"). - Discuss any interesting findings (like which signals were most important, or how early the model can predict a defect before end of cast). - Identify limitations of the current PoC. For instance, synthetic data might not capture all real-world variability; the model might need

retraining or adjustment when real data is available; sensors in reality can be noisy or fail, so robustness is key. - Suggest improvements (maybe integrating additional data like image analysis of the slab surface if available, or using more advanced algorithms if needed).

Throughout implementation, we will maintain good software practices: using version control (git) for code, structuring code into functions/modules (e.g., data_gen.py, train_model.py, infer_real_time.py), writing at least basic documentation for each, and possibly small unit tests for data processing functions. This not only makes the demo more production-like, it will impress interviewers with attention to code quality.

## Real-Time Integration and Deployment Considerations

Although the initial focus is the PoC, it's crucial to outline how this system would be deployed in a real production environment at a steel plant like those served by KISS Technologies. Below are considerations and recommendations for a production-quality system:

- **Industrial Connectivity:** The system should connect to existing plant data infrastructure. Likely, a continuous caster has a Level 1 (PLC) and Level 2 (process computer) system. We can tap into these via OPC UA servers or similar, or use IoT gateways to stream data. For a modern solution, using a message broker (MQTT or Kafka) is advisable for decoupling – sensors publish data which our ML system subscribes to. We should ensure low-latency, reliable data transfer. If using Kafka, for example, we would create topics for "caster_sensors" and possibly partition by sensor or cast ID.

- **Scalability and Microservices:** In production, it's better to break the system into microservices. For instance, one service handles data ingestion and preprocessing, another service hosts the ML model inference (as a REST API or gRPC service), and another might handle the dashboard or alerting. This allows each to be scaled or updated independently. Our PoC is likely a single-process script, but we can highlight that it would be straightforward to refactor into separate components, possibly containerized with Docker. Each container could be deployed on-premises or in the cloud depending on latency and data governance requirements (steel plants often require on-site processing for real-time control).

- **Model Deployment and Serving:** For the PyTorch model, options include embedding it in a Python service (using something like FastAPI or Flask to serve predictions) or using specialized servers like TorchServe. Given the need for real-time inference (on the order of 1-second or faster), a lightweight server or even in-process inference in a streaming job is fine. We should ensure the model is loaded once and reused for efficiency. If multiple caster machines or multiple strands are to be served, the system should handle concurrent inference (which is a matter of running the model on multiple threads or having a model per process). Also consider using a GPU if the model is large, but in our case CPU should suffice.

- **Data Frequency and Volume:** Continuous casting can generate a lot of data (multiple sensors at 1 Hz or faster for hours). Our system must handle this volume without bottleneck. Techniques include buffering data in memory for one cast then discarding after use, compressing or downsampling less critical signals, and ensuring any databases or logs can write data fast enough. We may implement a circular buffer for recent data if doing rolling predictions.

- **Fault Tolerance:** Production systems need to handle sensor failures, network hiccups, or model errors gracefully. We should incorporate basic fault tolerance: e.g., if a sensor reading is missing, use the last known value or a safe default; if the ML model service crashes, have an automatic restart (using a supervisor or container restart policy) and a fallback logic (perhaps a simpler rule-based alarm) so that operators are not left blind.

- **Monitoring and Maintenance:** Once deployed, the model's performance should be monitored. Over time, the process might change (different steel grades, equipment upgrades) leading to **model drift**. We should set up monitoring of the model's predictions vs actual outcomes (once actual quality is known from inspections) to recalibrate the model if needed. Including a feedback loop where false alarms or missed defects are analyzed will help improve the system. We can also log a **confidence metric** or anomaly scores to gauge when the model is unsure.

- **Security and Access:** In a plant setting, security is paramount. The data pipeline should be secure (encrypted connections from PLCs to our system, etc.). If any cloud component is used (for remote monitoring or storage), ensure compliance with company policies. For the PoC, security is less of an issue, but demonstrating awareness (e.g., not exposing a open port without auth in the demo) is good.

- **User Interface and Alerting:** We touched on visualization for the demo, but in production one might integrate with the plant's existing HMI or SCADA system. For example, the prediction could be sent to the operator's console or as an alarm in the control room if quality deviates. Additionally, storing predictive data to a database allows engineers to analyze trends (did a certain parameter drift over days and start causing defects?). KISS's own platform ("CasterAnalytics" [22] ) likely provides a UI for this; our system should be able to feed into such a platform or replicate key aspects (charts, reports).

By addressing these points, we ensure the solution is not just a one-off demo, but a blueprint for a real **production-quality deployment**. The modular design and choice of widely-used technologies (Python, PyTorch, message queues, etc.) make it feasible to harden and integrate this system in an industrial setting. We would also prepare documentation and handover materials as part of production readiness – e.g., a README for how to deploy the system, configuration files for different caster setups, and a maintenance guide for updating the model with new data.

## Conclusion

This technical specification has outlined a comprehensive plan for developing a **Predictive Quality Monitoring System** for continuous casting in a steel plant. By leveraging sensor data (temperature, speed, oscillation, etc.) and machine learning algorithms, the system can **predict quality issues in real-time**, enabling operators to take proactive measures to reduce defects [23] . We have covered the domain context of continuous casting, identified the relevant data sources (process parameters, material properties, environment, equipment status, and quality outcomes), and proposed a robust system architecture separating data ingestion, processing, modeling, and output visualization.

We evaluated different modeling approaches – from interpretable **tree-based models** to advanced **deep learning sequence models** – and recommended an implementation that uses both, balancing quick insights and powerful pattern recognition [24] [19] . The step-by-step plan ensures that we can build a

working PoC quickly (using synthetic data to simulate plant operations) and incrementally improve it. Emphasis was placed on making the demo impressive: real-time prediction updates, visual feedback, and a discussion of results will show a practical grasp of the problem and solution.

Finally, we addressed how this PoC can transition to a production system at KISS/steel plant: considerations like streaming data integration, microservice architecture, deployment, and ongoing maintenance were detailed to demonstrate foresight in engineering a **production-grade solution**. By following this spec, the resulting demo will not only function as a prototype but also serve as a conversation piece on how to implement cutting-edge **data-driven quality control** in steel manufacturing – aligning perfectly with KISS Technologies' mission of using data to take the guesswork out of casting and improve steel quality [2] [22] .

**References (from Research & Domain)**: This design is informed by industry knowledge and recent research. Continuous casting quality is influenced by many factors [2] [16] , and modern data-driven approaches have shown success in predicting defects like rhomboidity and inclusions using machine learning [15] [25] . KISS Technologies' own insights on the importance of data (process parameters, material, environment, equipment, production) guided our data collection strategy [26] [12] . Prior studies have used techniques from random forests to CNN-LSTM networks to tackle similar problems [24] [19] , which we leveraged in choosing our model approach. This fusion of domain research and technical planning should result in a compelling demonstration system for the interview.

---

[1] [6] [7] [13] [14] [17] [19] [20] [21] [24] [25] Machine-Learning Algorithms for Process Condition Data-Based Inclusion Prediction in Continuous-Casting Process: A Case Study
https://www.mdpi.com/1424-8220/23/15/6719

[2] [3] Kiss Technologies - Continuous Casting Process Technology
https://kisstechnologies.com/

[4] [5] [15] Enhancing Continuous Casting Efficiency with Machine Learning
https://econ-tech.com/blog-n1/enhancing-continuous-casting-efficiency-with-machine-learning

[8] [9] [10] [11] [12] [16] [18] [22] [23] [26] What data should a steelmaker use for a modern data-driven continuous casting machine? - Kiss Technologies
https://kisstechnologies.com/what-data-should-a-steelmaker-use-for-a-modern-data-driven-continuous-casting-machine/