

# VLE Modeling of Ethyl Acetate–Methanol via NRTL Parameter Estimation and Residual Multilayer Perceptron (MLP) Neural Network Correction

## Overview

The experimental **VLE data for the ethyl acetate–methanol system** were validated using the van Ness point-to-point consistency test. Group-contribution baselines (UNIFAC and ASOG) were computed to provide parameter-independent estimates of activity coefficients for comparison. In this study, the system was modeled with **NRTL** using literature parameters and refitted binary interaction parameters, alongside a **machine-learning residual multilayer perceptron (MLP)** that learns corrections to the NRTL+Raoult vapor-phase composition. Consequently, the models' accuracies were quantified against the experimental data using MAE, RMSE, and SMAPE.

## Objectives

1. To utilize NRTL model using literature binary interaction parameters  $g_{12}$  and  $g_{21}$  and generate a baseline VLE.
2. To estimate the NRTL binary interaction parameters and generate a VLE prediction using the fitted values.
3. To train a residual neural network (MLP) to adjust the NRTL baseline vapor-phase composition and generate a hybrid NRTL–NN prediction.
4. To compare Baseline, Fitted, and NRTL-NN predictions using MAE, RMSE, and SMAPE

## Libraries Import

- Imports core packages
- Sets global constants and plotting style (matplotlib only)

In [12]:

▶

```
1  import math, json, warnings
2  import numpy as np
3  import pandas as pd
4  import matplotlib.pyplot as plt
5  from thermo import Chemical
6
7  from dataclasses import dataclass
8  from typing import Tuple, Dict
9
10 from sklearn.neural_network import MLPRegressor
11 from sklearn.preprocessing import StandardScaler
12 from sklearn.pipeline import Pipeline
13 from sklearn.model_selection import KFold
14 from sklearn.metrics import mean_absolute_error, root_mean_squared_error
15
16 from scipy.optimize import minimize
17
18 R = 8.314462618 # J/mol/K
```

## ETL (Extract, Transform, and Load)

- Importing dataset which contains the experimental data
- Renames 'y2' to 'gamma2' because the paper's table column is  $\gamma_2$
- Adds constant pressure column  $P_{\text{kPa}} = 300$  (0.3 MPa dataset)

In [13]:

```
1 CSV_PATH = "ethyl_acetate_methanol_vle.csv" # change if your file is elsewhere
2
3 df = pd.read_csv(CSV_PATH)
4
5 # Renaming column
6 if "y2" in df.columns and "gamma2" not in df.columns:
7     df = df.rename(columns={"y2": "gamma2"})
8
9 # Safety checks
10 expected = {"T_K", "x1", "y1"}
11 missing = expected - set(df.columns)
12 if missing:
13     raise ValueError(f"Missing required columns: {missing}")
14
15 # Pressure (0.3 MPa) for this dataset
16 if "P_kPa" not in df.columns:
17     df["P_kPa"] = 300.0
18
19 # Derived
20 df["x2"] = 1.0 - df["x1"]
21 df.head()
```

Out[13]:

	T_K	x1	y1	gamma2	P_kPa	x2
0	369.45	0.000	0.000	1.000	300.0	1.000
1	368.85	0.022	0.035	0.989	300.0	0.978
2	368.55	0.031	0.049	0.993	300.0	0.969
3	368.15	0.044	0.063	1.005	300.0	0.956
4	367.65	0.077	0.102	1.014	300.0	0.923

## Psat(T): Antoine constants

- Uses the thermo package to get Psat(T) directly (in kPa)
- Component IDs: "ethyl acetate" and "methanol"

In [66]:

```
1 def Psat_kPa_component1(T_K: float) -> float:
2     chem = Chemical("ethyl acetate", T=T_K)
3     return chem.Psat/1000.0 # Pa → kPa
4
5 def Psat_kPa_component2(T_K: float) -> float:
6     chem = Chemical("methanol", T=T_K)
7     return chem.Psat/1000.0 # Pa → kPa
```

## NRTL (Baseline)

- Implements binary NRTL for γ1, γ2
- Uses data from Susial et al. (2012) where: α12=0.47, g12=1124.6 J/mol, g21=2202.2 J/mol

```
In [79]: 1 @dataclass
2 class NRTLParams:
3     alpha12: float
4     alpha21: float
5     g12: 1124.6      # J/mol
6     g21: 2202.2      # J/mol
7
8     # Symmetric non-randomness (alpha12 = alpha21 = 0.47):
9     NRTL_DEFAULT = NRTLParams(alpha12=0.47, alpha21=0.47, g12=1124.6, g21=2202.2)
10
11 def nrtl_gammas(x1: float, T_K: float, p: NRTLParams) -> Tuple[float, float]:
12     """
13     Binary NRTL activity coefficients (standard closed form).
14     Matches: ln G12 = -alpha12*tau12, ln G21 = -alpha21*tau21,
15             tau12 = g12/(RT), tau21 = g21/(RT)
16     """
17     # Guard against endpoints to avoid division-by-zero
18     x1 = float(np.clip(x1, 1e-12, 1.0 - 1e-12))
19     x2 = 1.0 - x1
20
21     tau12 = p.g12 / (R * T_K)
22     tau21 = p.g21 / (R * T_K)
23     G12 = math.exp(-p.alpha12 * tau12)
24     G21 = math.exp(-p.alpha21 * tau21)
25
26     S1 = x1 + x2 * G21
27     S2 = x2 + x1 * G12
28
29     # Textbook binary NRTL
30     ln_gamma1 = (x2**2) * (tau21 * (G21 / S1)**2 + (tau12 * G12) / (S2**2))
31     ln_gamma2 = (x1**2) * (tau12 * (G12 / S2)**2 + (tau21 * G21) / (S1**2))
32
33     return math.exp(ln_gamma1), math.exp(ln_gamma2)
```

## Modified Raoult's law

- Computes  $y_1$  from  $x_1$ ,  $T$ ,  $P$  using  $\gamma$  from NRTL and  $P_{\text{sat}}(T)$  from Antoine.

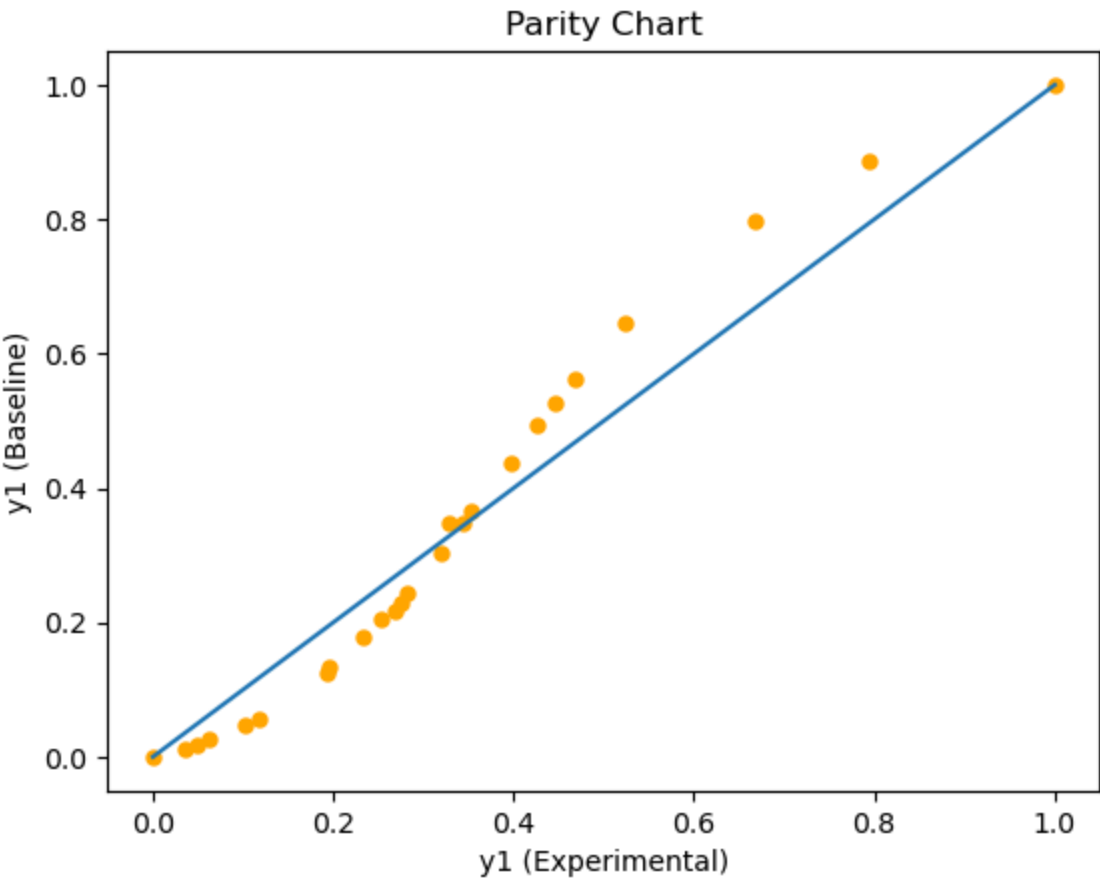
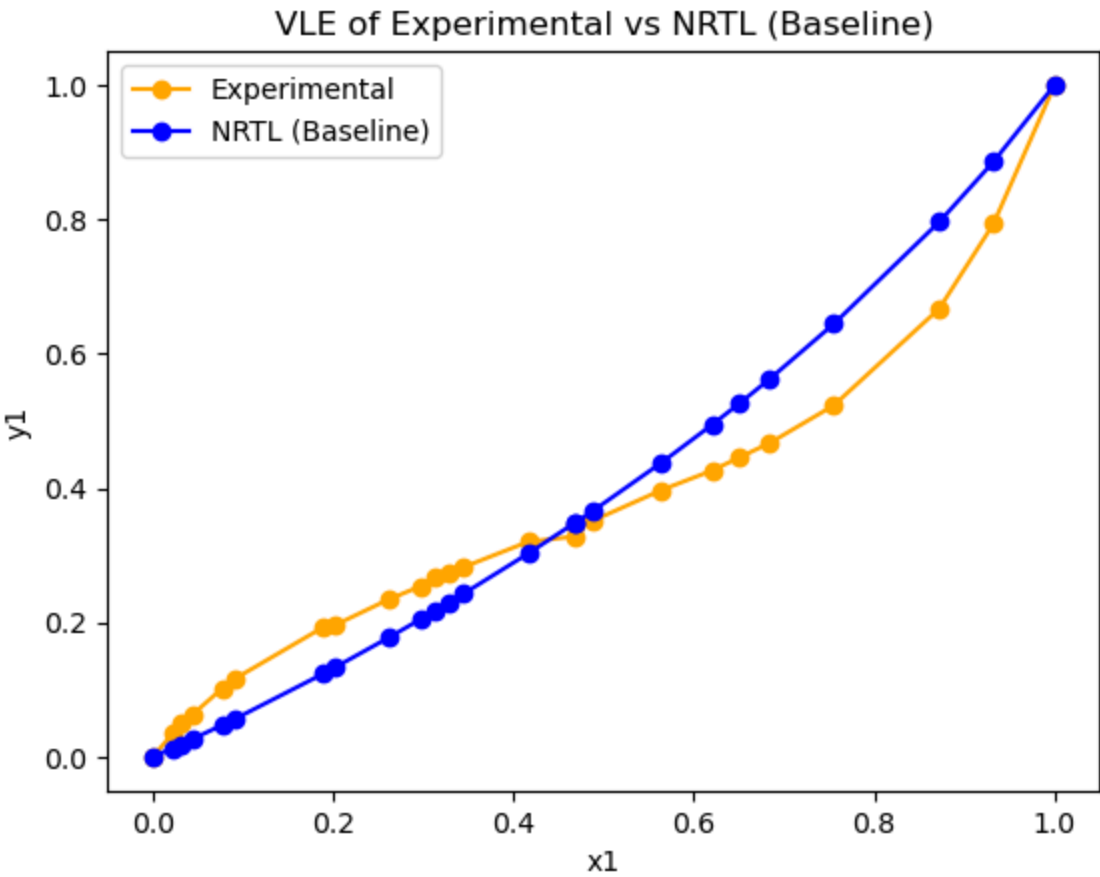
```
In [80]: 1 def predict_y1_row(x1: float, T_K: float, P_kPa: float, nrtl: NRTLParams) -> float:
2     x2 = 1.0 - x1
3     g1, g2 = nrtl_gammas(x1, T_K, nrtl)
4     P1 = Psat_kPa_component1(T_K)
5     P2 = Psat_kPa_component2(T_K)
6     num1 = x1 * g1 * P1
7     num2 = x2 * g2 * P2
8     y1 = num1 / (num1 + num2 + 1e-16)
9     return max(0.0, min(1.0, y1))
```

## Baseline predictions and metrics/plots

- Vectorizes  $y_1$  prediction over the dataframe.
- Computes MAE/RMSE/AARD%.
- Plots  $y$ - $x$  and parity plots.

```
In [143]: 1 def eval_baseline(df_in: pd.DataFrame, nrtl: NRTLParams):
2     df = df_in.copy()
3     df["y1_calc_nrtl"] = [
4         predict_y1_row(x, T, P, nrtl)
5         for T, x, P in zip(df["T_K"], df["x1"], df["P_kPa"])
6     ]
7
8     y_true = df["y1"].to_numpy()
9     y_pred = df["y1_calc_nrtl"].to_numpy()
10
11     mae = mean_absolute_error(y_true, y_pred)
12     rmse = root_mean_squared_error(y_true, y_pred)
13     aard = 100.0 * np.mean(np.abs((y_true - y_pred) / (y_true + 1e-12)))
```

```
In [144]: 1 # Plots
2 plt.figure()
3 plt.plot(df["x1"], df["y1"], label="Experimental", color="orange", marker='o',
4 plt.plot(df_pred["x1"], df_pred["y1_calc_nrtl"], label="NRTL (Baseline)", color
5 plt.xlabel("x1"); plt.ylabel("y1")
6 plt.title("VLE of Experimental vs NRTL (Baseline)")
7 plt.legend()
8 plt.show()
9
10 plt.figure()
11 plt.scatter(df_pred["y1"], df_pred["y1_calc_nrtl"], s=25, color="orange")
12 plt.plot([0,1],[0,1])
13 plt.xlabel("y1 (Experimental)"); plt.ylabel("y1 (Baseline)")
14 plt.title("Parity Chart")
15 plt.show()
```



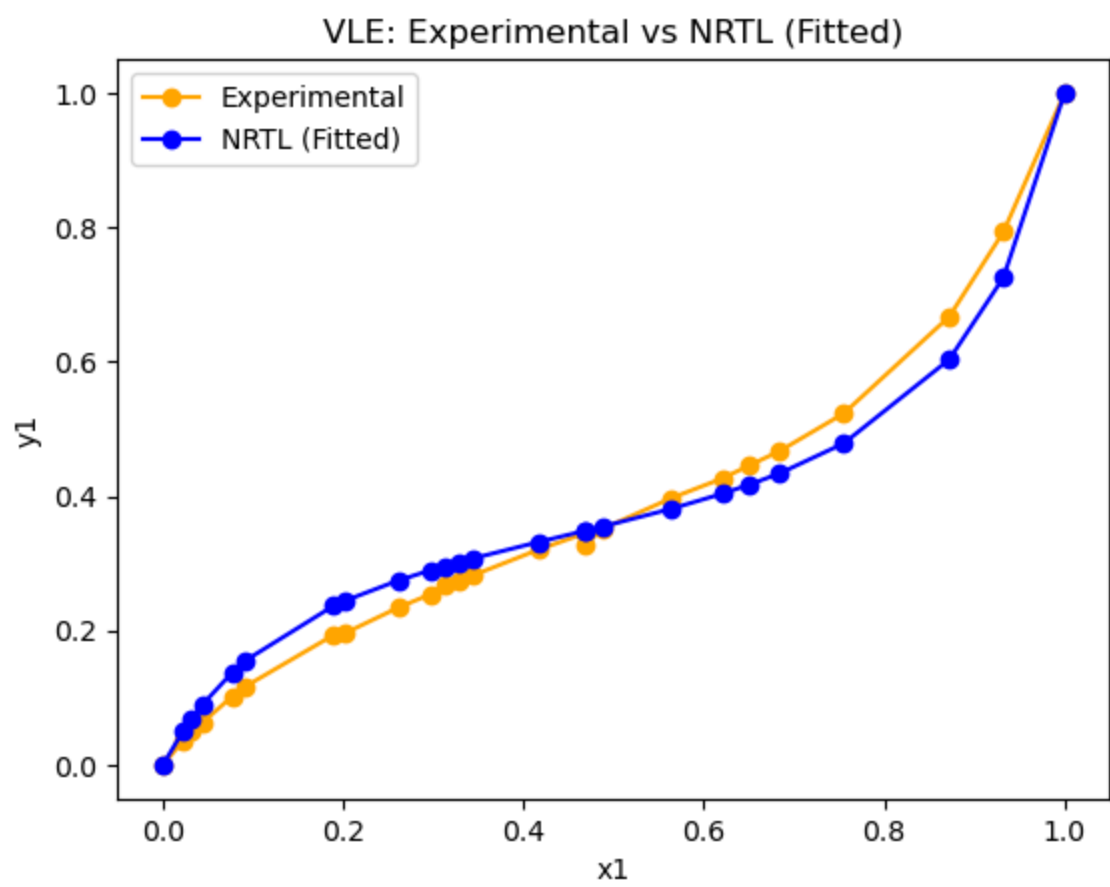
## NRTL (Fitted)

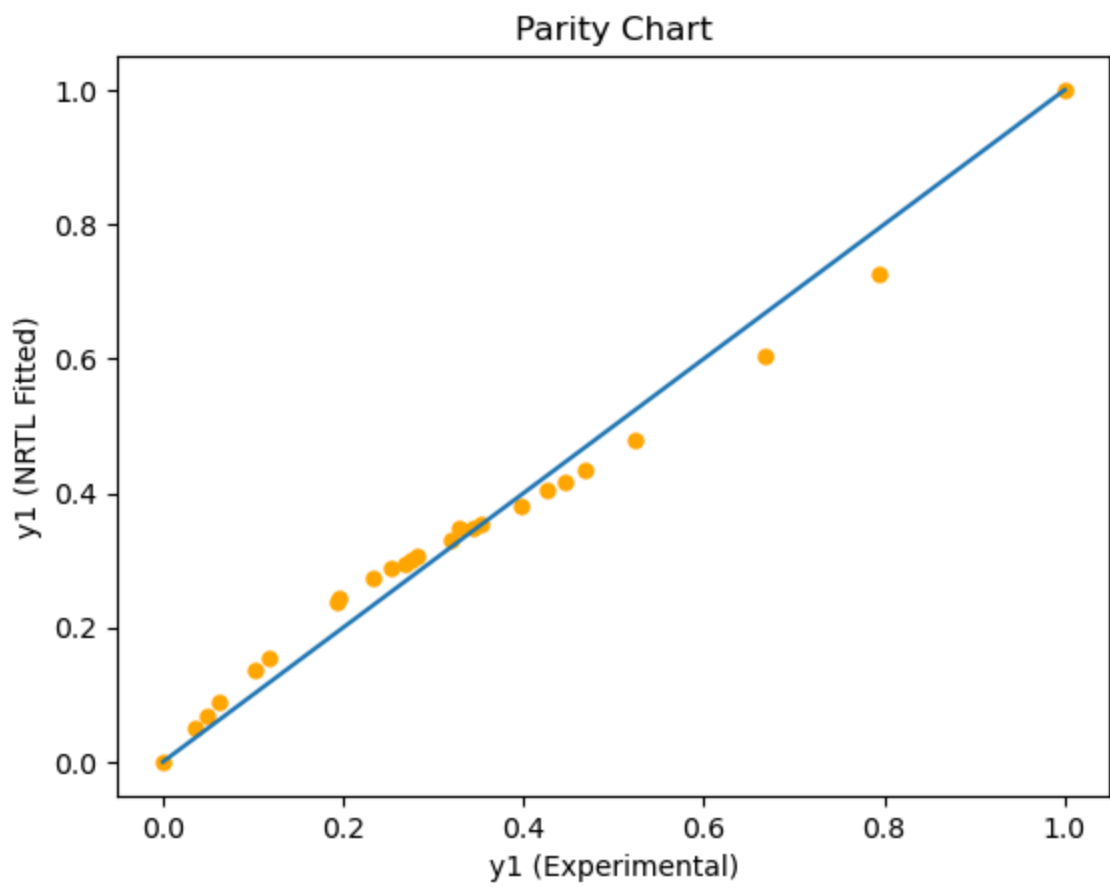
- Estimates new values for  $g_{12}$ ,  $g_{21}$  but is based on  $\alpha = 0.47$  (Susial et al., 2012) to minimize SSE on  $y_1$ .
- Prints fitted params and new metrics.

In [131]: ▶

```
1  # Description:
2  # - Fits g12, g21 (alpha fixed) via L-BFGS-B on SSE of y1
3
4  def sse_for_params(theta: np.ndarray, df_local: pd.DataFrame, alpha: float) ->
5      g12, g21 = theta
6      p = NRTLParams(alpha=alpha, g12=g12, g21=g21)
7      yhat = [predict_y1_row(x, T, P, p) for T, x, P in zip(df_local["T_K"], df_lo
8      y = df_local["y1"].to_numpy()
9      return float(np.sum((y - np.array(yhat))**2))
10
11 x0 = [2000, 3000]
12 bnds = [(100, 15000), (100, 15000)]
```

```
1 # Plot
2
3 # Sort for smooth lines
4 order = np.argsort(df["x1"].to_numpy())
5 x_sorted = df["x1"].to_numpy()[order]
6 yexp_sorted= df["y1"].to_numpy()[order]
7 yfit_sorted= y_fit[order]
8
9 # Experimental (line+dots) vs NRTL (Fitted)
10 plt.figure()
11 plt.plot(x_sorted, yexp_sorted, label="Experimental", color="orange", marker='o')
12 plt.plot(x_sorted, yfit_sorted, label="NRTL (Fitted)", color="blue", marker='o')
13 plt.xlabel("x1"); plt.ylabel("y1")
14 plt.title("VLE: Experimental vs NRTL (Fitted)")
15 plt.legend(); plt.show()
16
17 # Parity plot using y_fit
18 plt.figure()
19 plt.scatter(df["y1"], y_fit, s=25, color="orange")
20 plt.plot([0,1],[0,1])
21 plt.xlabel("y1 (Experimental)"); plt.ylabel("y1 (NRTL Fitted)")
22 plt.title("Parity Chart")
23 plt.show()
```





## NRTL (Residual MLP)

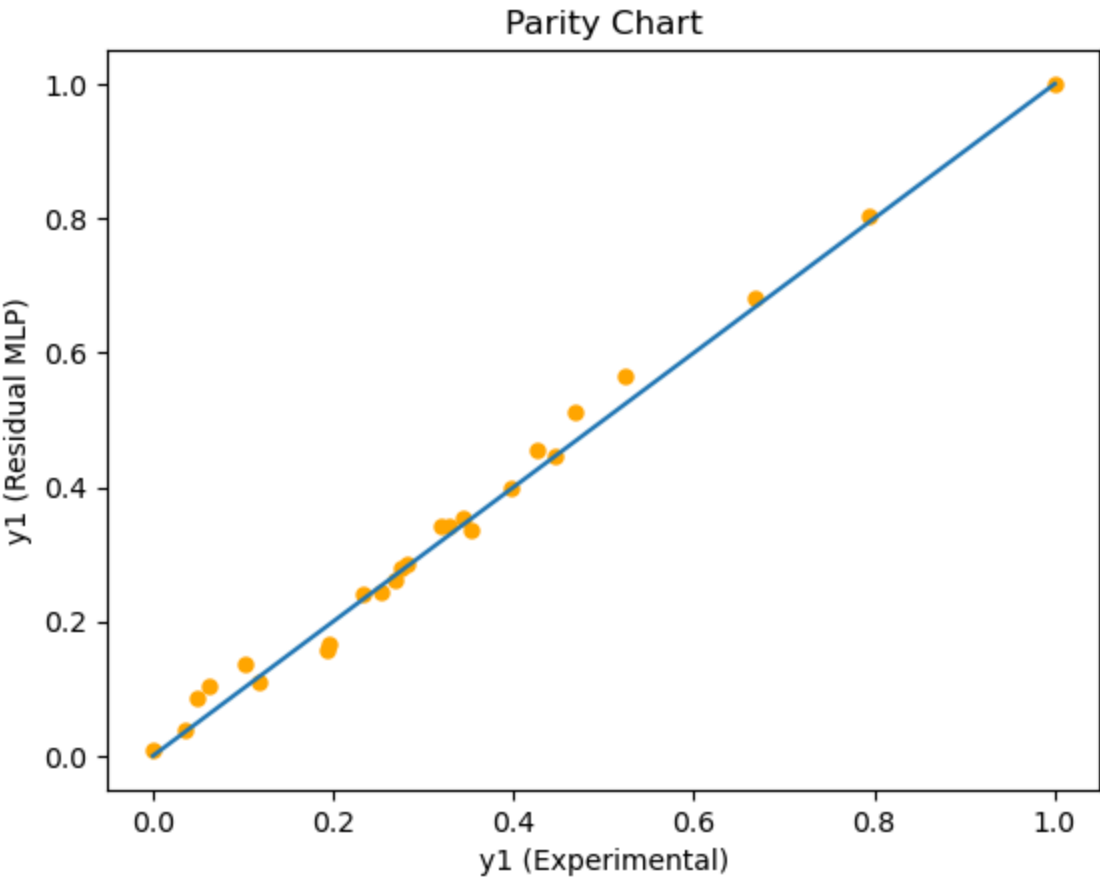
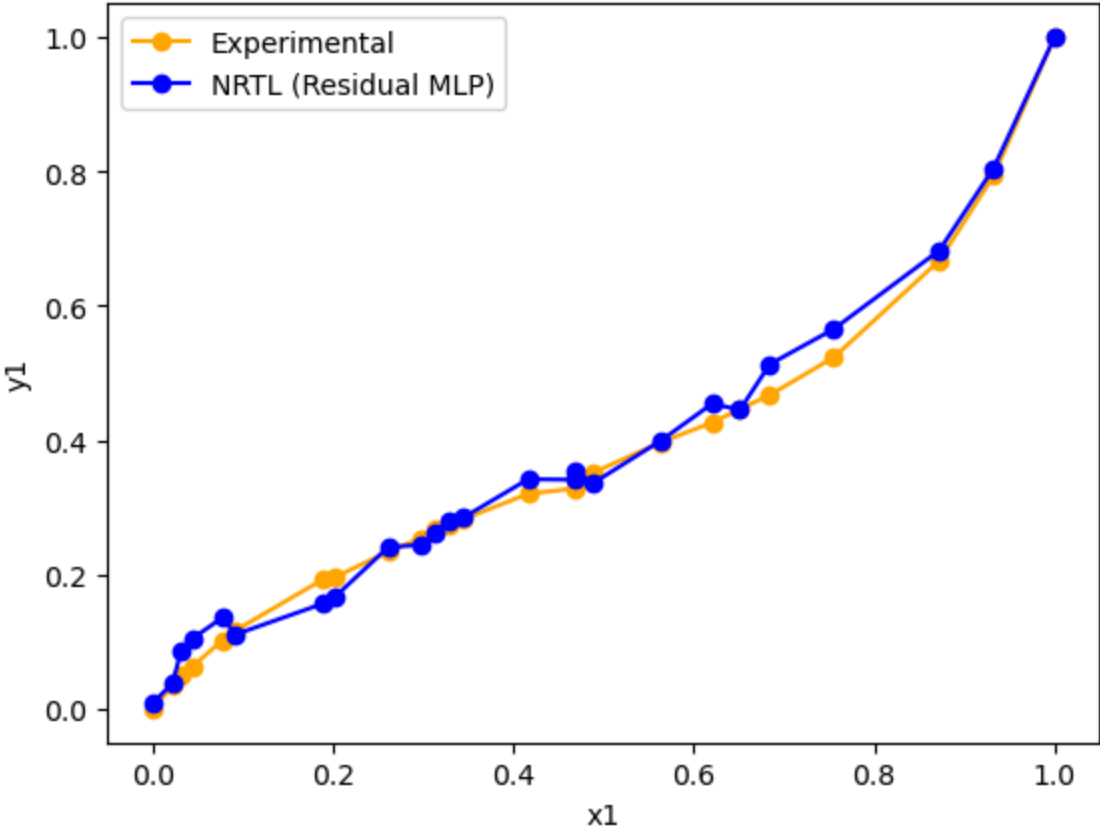
- Learn residual  $\Delta y = y_{\text{exp}} - y_0$  where  $y_0 = \text{NRTL} + \text{Raoult prediction}$
- Inputs:  $[x_1, T_{\text{K}}, P_{\text{kPa}}, y_0, \text{gamma1\_NRTL}, \text{gamma2\_NRTL}]$
- Model: scikit-learn MLPRegressor
- Evaluation: 5-fold CV; report MAE, RMSE, AARD% and plot vs baseline

In [84]:

```
1 def build_residual_features(df_in, nrtl_params):
2     df0 = df_in.copy()
3     # baseline y0 and NRTL gammas
4     y0, g1_list, g2_list = [], [], []
5     for T, x, P in zip(df0["T_K"], df0["x1"], df0["P_kPa"]):
6         y0_i = predict_y1_row(x, T, P, nrtl_params)
7         g1_i, g2_i = nrtl_gammas(x, T, nrtl_params)
8         y0.append(y0_i); g1_list.append(g1_i); g2_list.append(g2_i)
9     y0 = np.array(y0, dtype=np.float32)
10    g1 = np.array(g1_list, dtype=np.float32)
11    g2 = np.array(g2_list, dtype=np.float32)
12
13    X = np.column_stack([
14        df0["x1"].to_numpy(dtype=np.float32),
15        df0["T_K"].to_numpy(dtype=np.float32),
16        df0["P_kPa"].to_numpy(dtype=np.float32),
17        y0, g1, g2
18    ])
19    y = df0["y1"].to_numpy(dtype=np.float32)
20    return X, y, y0
21
22 # Build features/targets
23 X, y, y0 = build_residual_features(df, NRTL_DEFAULT)
24
25 # 5-fold CV fit of residuals
26 kf = KFold(n_splits=5, shuffle=True, random_state=42)
27 yhat_nn = np.zeros_like(y, dtype=np.float32)
28
29 for tr, te in kf.split(X):
30     model = Pipeline(steps=[
31         ("scaler", StandardScaler()),
32         ("mlp", MLPRegressor(
33             hidden_layer_sizes=(64, 64),
34             activation="relu",
35             alpha=1e-4, # L2 regularization
36             learning_rate_init=1e-3,
37             max_iter=5000,
38             early_stopping=True,
39             n_iter_no_change=100,
40             random_state=42
41         ))
42     ])
43     # Train on residuals
44     model.fit(X[tr], (y[tr] - y0[tr]))
45     # Predict residuals and add back baseline
46     res_pred = model.predict(X[te])
47     yhat_nn[te] = np.clip(y0[te] + res_pred.astype(np.float32), 0.0, 1.0)
48
49 # Metrics vs experimental
50 mae = mean_absolute_error(y, yhat_nn)
51 rmse = root_mean_squared_error(y, yhat_nn)
52 aard = 100.0 * np.mean(np.abs((y - yhat_nn) / (y + 1e-12)))
```



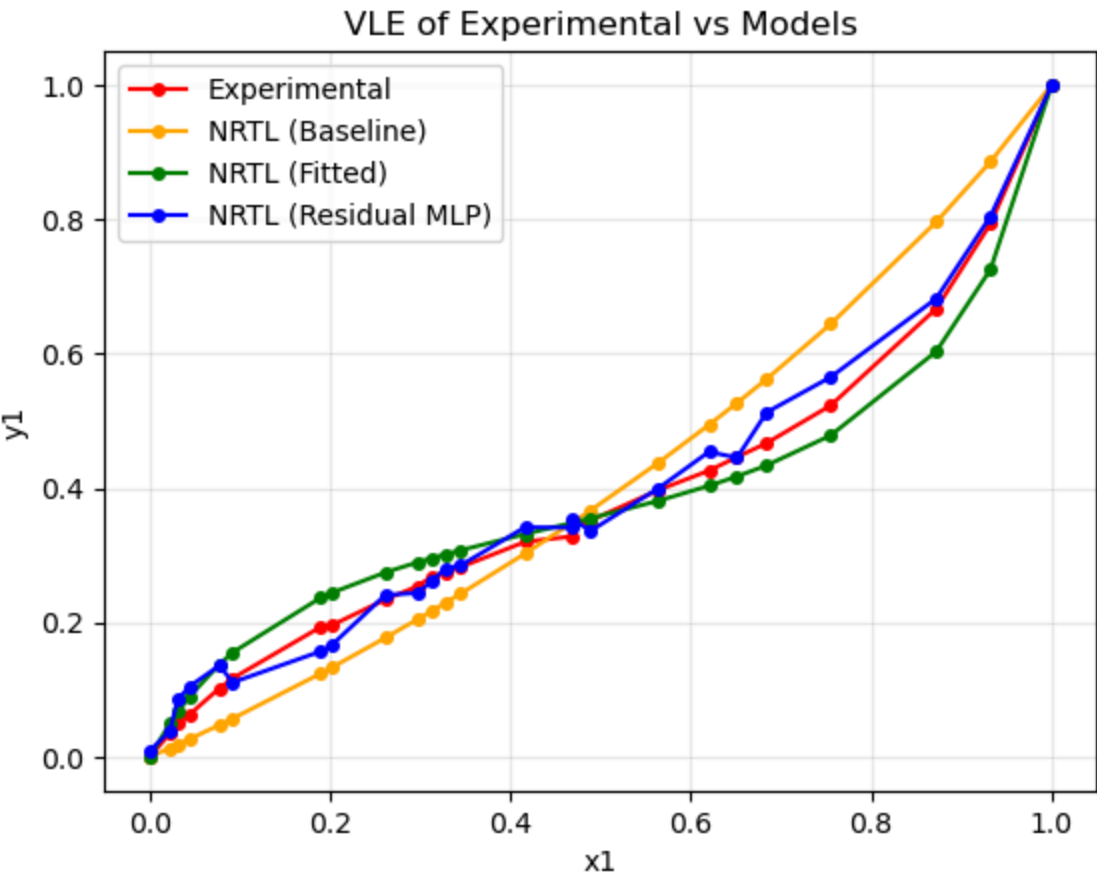
```
In [102]: 1 # Plot exp vs baseline vs residual-MLP
2 plt.figure()
3 plt.plot(df["x1"], df["y1"], label="Experimental", color="orange", marker='o',
4 plt.plot(df["x1"], yhat_nn, label="NRTL (Residual MLP)", color="blue", marker='
5 plt.xlabel("x1"); plt.ylabel("y1")
6 plt.legend()
7 plt.show()
8
9 # Parity plot
10 plt.figure()
11 plt.scatter(y, yhat_nn, s=25, color="orange")
12 plt.plot([0,1],[0,1])
13 plt.xlabel("y1 (Experimental)"); plt.ylabel("y1 (Residual MLP)")
14 plt.title("Parity Chart")
15 plt.show()
```



Summary of Results

In [142]:

```
1 # Plot
2 plt.plot(df["x1"], df["y1"], label="Experimental", color="red", marker='o', mar
3 plt.plot(df_pred["x1"], df_pred["y1_calc_nrtl"], label="NRTL (Baseline)", color
4 plt.plot(x_sorted, yfit_sorted, label="NRTL (Fitted)", color="green", marker='o
5 plt.plot(df["x1"], yhat_nn, label="NRTL (Residual MLP)", color="blue", marker='
6
7 plt.xlabel("x1")
8 plt.ylabel("y1")
9 plt.title("VLE of Experimental vs Models")
10 plt.legend()
11 plt.grid(True, alpha=0.3)
12 plt.show()
```



In [149]:

```
1 def smape_percent(y_true, y_pred, epsilon=1e-10):
2     denominator = np.abs(y_true) + np.abs(y_pred) + epsilon
3     return 100 * np.mean(
4         2 * np.abs(y_pred - y_true) / denominator
5     )
6
7 def compute_metrics(y_true, y_pred):
8     return {
9         "MAE": mean_absolute_error(y_true, y_pred),
10        "RMSE": root_mean_squared_error(y_true, y_pred),
11        "SMAPE%": smape_percent(y_true, y_pred),
12        "SSE": np.sum((y_true - y_pred) ** 2),
13    }
14
15 y_exp = df["y1"].to_numpy()
16
17 # NRTL (Baseline)
18 try:
19     y_base = df_pred["y1_calc_nrtl"].to_numpy()
20 except Exception:
21     y_base = np.array([
22         predict_y1_row(x, T, P, NRTL_DEFAULT)
23         for T, x, P in zip(df["T_K"], df["x1"], df["P_kPa"])
24     ])
25
26 rows = [("Baseline (NRTL)", compute_metrics(y_exp, y_base))]
27
28 # NRTL (Fitted)
29 try:
30     p_fit = NRTLParams(
31         alpha12=NRTL_DEFAULT.alpha12,
32         alpha21=NRTL_DEFAULT.alpha21,
33         g12=g12_fit,
34         g21=g21_fit
35     )
36     y_fit = np.array([
37         predict_y1_row(x, T, P, p_fit)
38         for T, x, P in zip(df["T_K"], df["x1"], df["P_kPa"])
39     ])
40     rows.append(("NRTL (Fitted)", compute_metrics(y_exp, y_fit)))
41 except NameError:
42     pass # g12_fit/g21_fit not defined yet
43
44 # NRTL (Residual MLP)
45 if "yhat_nn" in globals():
46     rows.append(("NRTL (Residual MLP)", compute_metrics(y_exp, yhat_nn)))
47
48 # Summary DataFrame
49 summary_df = pd.DataFrame(
50     [(name, m["MAE"], m["RMSE"], m["SMAPE%"], m["SSE"]) for name, m in rows],
51     columns=["Model", "MAE", "RMSE", "SMAPE%", "SSE"]
52 )
53
54 display(
55     summary_df.style.format({
56         "MAE": "{:.4f}",
57         "RMSE": "{:.4f}",
58         "SMAPE%": "{:.2f}",
59         "SSE": "{:.4f}"
60     })
61 )
```

	Model	MAE	RMSE	SMAPE%	SSE
0	Baseline (NRTL)	0.0505	0.0609	28.66	0.0927
1	NRTL (Fitted)	0.0278	0.0330	12.80	0.0272
2	NRTL (Residual MLP)	0.0173	0.0226	17.65	0.0127

Interpretation of Results

Experimental vs. models:

- The experimental points form a smooth, convex VLE curve. All three models reproduce the overall shape and the correct limits at very low and very high x1

**Baseline NRTL (g12 and g21 specified)**

- Captures the trend but tends to overpredict y1 at medium–high x1

**Fitted NRTL (re-estimated g12 and g21 with  $\alpha$  fixed)**

- Moves the curve slightly closer to the data; the data improved significantly but can't still capture the experimental data

**NRTL + Residual Neural Network (MLP)**

- Learns the remaining error of the NRTL prediction and adds a correction. This hybrid follows the experimental points more closely across most compositions.

**Implications**

- Baseline NRTL provides a sound thermodynamic reference.
- Fitting reduces bias but only slightly.
- The hybrid NRTL + NN delivers the best accuracy for interpolation within the data's temperature/condition range. But, it should not be extrapolated beyond that range without additional data or retraining.

In [150]:



```
1  # Save Modeling Prediction in CSV File
2
3  # Create a DataFrame with all results
4  results_df = pd.DataFrame({
5      "x1_exp": df["x1"],
6      "y1_exp": df["y1"],
7      "x1_baseline": df_pred["x1"],
8      "y1_baseline": df_pred["y1_calc_nrtl"],
9      "x1_fitted": x_sorted,
10     "y1_fitted": yfit_sorted,
11     "x1_nn": df["x1"],
12     "y1_nn": yhat_nn
13 })
14
15 # Save to CSV
16 results_df.to_csv("[G03] Model Results.csv", index=False)
17
18 print("Results saved to [G03] Model Results.csv")
```

Results saved to [G03] Model Results.csv