



# Machine Learning Inference Tuning using Quantum XL Version 18

Ray Balisnomo

Friday, April 17, 2026

# Introduction

- Machine Learning (ML) inference tuning is the systematic optimization of how a trained model executes in production to achieve predictable, robust performance while respecting accuracy, resource, and reliability constraints.
- ML inference tuning is the practice of deliberately adjusting how a trained machine learning model is executed (not how it is trained) to achieve the best trade off between performance, accuracy, resource usage, and robustness in a real deployment.
- Think of it as engineering optimization, not data science.
- Inferencing tuning asks: “What settings work **consistently**, not just once?”

# ML Inference Tuning vs. “Benchmark Chasing”

## **BENCHMARK CHASING**

- Maximize one metric
- Best-case scenario
- Fixed workload
- Hero numbers (works just once)

## **INFERENCE TUNING**

- Balance multiple objectives
- Robust across conditions
- Workload variability
- Predictable performance (works consistently)

# Why engineers increasingly use DOE for inference tuning

- 1) Because inference behavior is:
  - **nonlinear**
  - **hardware dependent**
  - **interaction heavy**
- 2) Response Surface Design allows you to:
  - **quantify trade offs**
  - **predict unseen settings**
  - **find robust operating regions**
  - **defend decisions with data**
- 3) It turns ML inference from **trial and error** into **engineering discipline**.

# Problem Context

You are deploying an ML model (CNN / transformer / classical ML) in a firmware or software runtime where you must balance:

- **inference latency**
- **throughput**
- **accuracy**
- **CPU / accelerator utilization**
- **power consumption**

These tradeoffs are nonlinear and highly interactive, so OFAT tuning (one factor at a time) fails.

## Factors and Engineering Ranges

- The ranges below are chosen to reflect what engineers can actually change without redesigning the model.
- Control factors: Response Surface Design (RSD) does **not** allow categorical factors.

<b>Factor</b>	<b>Symbol</b>	<b>Type</b>	<b>Low</b>	<b>Center</b>	<b>High</b>
<b>Batch Size</b>	<b>A</b>	<b>Numeric</b>	<b>1</b>	<b>8</b>	<b>16</b>
<b>Quantization</b>	<b>B</b>	<b>Ordinal</b>	<b>INT4</b>	<b>INT8</b>	<b>FP16</b>
<b>Precision Mode</b>	<b>C</b>	<b>Ordinal</b>	<b>FP16</b>	<b>TF32</b>	<b>FP32</b>

# Encoding Ordinal Factors

- Although quantization and precision are categorical, they are **ordered by numeric fidelity** and can be treated as **ordinal numeric factors** for RSD.
- This is common practice when strategies represent increasing strength or fidelity.

Level	Quantization	Precision Mode
-1	INT4	FP16
0	INT8	TF32
+1	FP16	FP32

## Design Type: Face-Centered Central Composite Design (CCD)

- Three factors
- Curvature expected (batch size saturation)
- Strong interactions (precision mode x batch)
- No unsafe extrapolation beyond feasible runtime settings
- Total runs = 20
  - Runs 1-8 are full factorials
  - Runs 9-14 are face-centered axial points
  - Runs 15-20 are center-point replications to estimate noise and drift.

Design Wizard | **Create Design** | Modify Design | Analyze Design | Optimize | Charts | ProTest (Covering Array) | Options | About

# QXL DOE > Create Design > Create CCD Design

- Modeling**
  - Create 2-Level Factorial Design
  - Create 3-Level Factorial Design
  - Create N-Level Factorial Design
  - Create CCD Design**
  - Create Box-Behnken Design
- Screening**
  - Create Plackett-Burman Design
  - Create Taguchi Design
- Special**
  - Setup Historical Analysis
  - Create Custom Design
  - Create D-Optimal Design
- Import**
  - Import From DOE PRO XL

**Create CCD Design**  
Create blocked or unblocked Central Composite Design.

Q7  
A  
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19

# Choose Full factorial, 3 Factors in 8 Runs. Resolution: Full

Central Composite Design

Select the base design

Recommended Designs Use with caution

		Factors									
		2	3	4	5	6	7	8	9	10	
factorial	Full	Full	Full	Full	Full	Full	Full				
	Half		III	IV	V	VI	VII	VIII			
	Quarter				IV	IV	IV	V	VI		
	Eighth					III	IV	IV	IV	IV	V
	Sixteenth						III	IV	IV	IV	IV

Full factorial, 3 Factors in 8 Runs.  
Resolution: Full.

Help Cancel Next > Finish

Central Composite Design

### Enter Factor Names and Coding

For Quantitative Inputs enter the low and high values (e.g. low=10 high=20)

Name	Low	High
A	-1	1
B	-1	1
C	-1	1

Help    Cancel    < Back    Next >    Finish

Coded units

Coded units

Coded units

Central Composite Design

### Enter Factor Names and Coding

For Quantitative Inputs enter the low and high values (e.g. low=10 high=20)

Name	Low	High
Batch Size Rows	1	16
Quantization	-1	1
Precision Mode	-1	1

Level	Quantization	Precision Mode
-1	INT4	FP16
0	INT8	TF32
+1	FP16	FP32

Buttons: Help, Cancel, < Back, Next >, Finish

Uncoded Units

Coded units

Coded units

Central Composite Design ×

### Define outputs (responses)

Select number of outputs:  ▾

Name: <input type="text" value="Output 1"/>	Type: <input type="text" value="Quantitative"/> ▾	<input type="checkbox"/> Has weights column
---	---	---

Design format

Table

Stacked

# Response Variables

Measure using repeated inference runs under identical input conditions.

<b>Response</b>	<b>Measurement Goal</b>	<b>Specification Limits</b>
<b>Y1 = Inference Latency</b>	<b>Minimize (milliseconds)</b>	<b>USL = 11 ms</b>
<b>Y2 = Throughput</b>	<b>Maximize (inferences per sec)</b>	<b>LSL = 125 ips</b>
<b>Y3 = Accuracy Drop</b>	<b>Minimize (percent)</b>	<b>USL = 2%</b>
<b>Y4 = Power or CPU Utilization</b>	<b>Constraint / Minimize</b>	<b>USL = 68%</b>

Central Composite Design

### Define outputs (responses)

Select number of outputs:

Name: <input type="text" value="Inference Latency ms"/>	Type: <input type="text" value="Quantitative"/>	<input type="checkbox"/> Has weights column
Name: <input type="text" value="Throughput inf/sec"/>	Type: <input type="text" value="Quantitative"/>	<input type="checkbox"/> Has weights column
Name: <input type="text" value="% Accuracy Drop"/>	Type: <input type="text" value="Quantitative"/>	<input type="checkbox"/> Has weights column
Name: <input type="text" value="% CPU Utilization"/>	Type: <input type="text" value="Quantitative"/>	<input type="checkbox"/> Has weights column

Design format

Table

Stacked

Help Cancel < Back Next > Finish

Central Composite Design

### Enter number of replicates and number of blocks

Replicates

Select number of replicates:

Select number of blocks:

Help Cancel < Back Next > Finish

Central Composite Design

### Axial and center points

Enter the number of center points:

Total number of center points: 6  
Total number of runs: 20

Alpha value

Manual

Optimal orthogonality (1.5246)

Rotatable (1.6818)

Face CCD Alpha = 1

Help Cancel < Back Finish

**Batch Size should be an integer; therefore, round down 8.5 to 8.**

	A	B	C
Run	Batch Size Rows	Quantization	Precision Mode
1	1	-1	-1
2	1	-1	1
3	1	1	-1
4	1	1	1
5	16	-1	-1
6	16	-1	1
7	16	1	-1
8	16	1	1
9	1	0	0
10	16	0	0
11	8.5	-1	0
12	8.5	1	0
13	8.5	0	-1
14	8.5	0	1
15	8.5	0	0
16	8.5	0	0
17	8.5	0	0
18	8.5	0	0
19	8.5	0	0
20	8.5	0	0

**Before**

	A	B	C
Run	Batch Size Rows	Quantization	Precision Mode
1	1	-1	-1
2	1	-1	1
3	1	1	-1
4	1	1	1
5	16	-1	-1
6	16	-1	1
7	16	1	-1
8	16	1	1
9	1	0	0
10	16	0	0
11	8	-1	0
12	8	1	0
13	8	0	-1
14	8	0	1
15	8	0	0
16	8	0	0
17	8	0	0
18	8	0	0
19	8	0	0
20	8	0	0

**After**



Run experiments and enter response variables.

Inference Latency ms		Throughput inf/sec		% Accuracy Drop		% CPU Utilization	
USL	<b>11</b>	USL		USL	<b>2</b>	USL	<b>68</b>
LSL		LSL	<b>125</b>	LSL		LSL	
<b>Data</b>		<b>Data</b>		<b>Data</b>		<b>Data</b>	
12		80		2.8		35	
18		60		0.3		55	
25		40		0.6		42	
30		35		0.1		70	
7		180		3.2		68	
11		140		0.5		88	
15		95		0.9		75	
20		85		0.2		95	
14		70		1.5		48	
8		165		1.9		82	
10		90		2.1		50	
17		65		0.7		60	
20		60		1.1		55	
16		75		0.4		72	
12		110		1.2		60	
13		108		1.1		62	
11		112		1.3		59	
12		109		1.2		61	
12		111		1.1		60	
13		107		1.2		63	

# Quantum XL

Central Composite Design

Base design: Full factorial

3 Factors in 20 Runs

6 Center points. Alpha = 1

Design is not replicated.

	A	B	C
Run	Batch Size Rows	Quantization	Precision Mode
1	1	-1	-1
2	1	-1	1
3	1	1	-1
4	1	1	1
5	16	-1	-1
6	16	-1	1
7	16	1	-1
8	16	1	1
9	1	0	0
10	16	0	0
11	8	-1	0
12	8	1	0
13	8	0	-1
14	8	0	1
15	8	0	0
16	8	0	0
17	8	0	0
18	8	0	0
19	8	0	0
20	8	0	0

edit response name (optional) -->

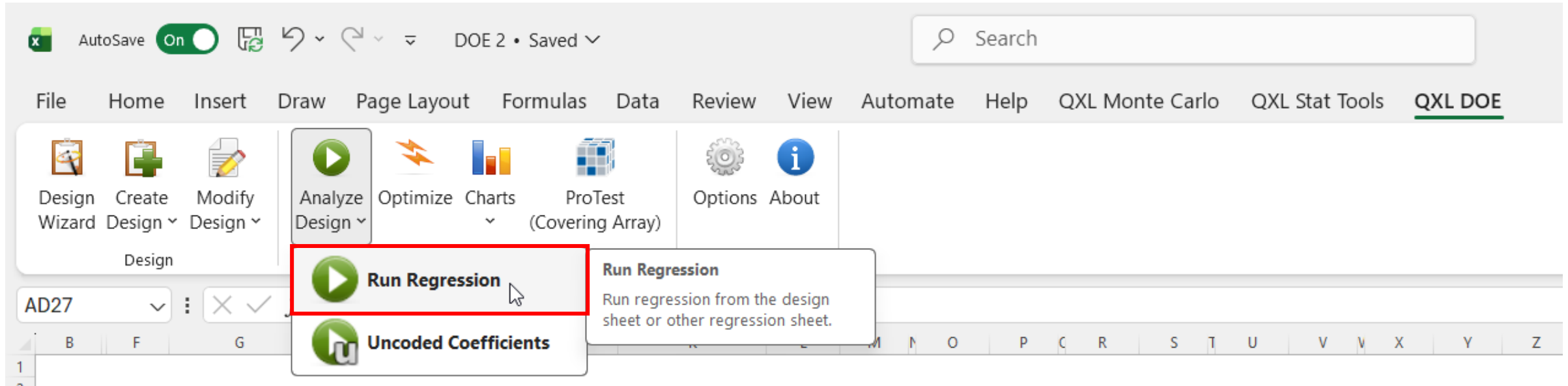
enter spec limits (optional) -->

<-- enter factor names

enter response data here -->

Inference Latency ms		Throughput inf/sec		% Accuracy Drop		% CPU Utilization	
USL	11	USL		USL	2	USL	68
LSL		LSL	125	LSL		LSL	
Data		Data		Data		Data	
12		80		2.8		35	
18		60		0.3		55	
25		40		0.6		42	
30		35		0.1		70	
7		180		3.2		68	
11		140		0.5		88	
15		95		0.9		75	
20		85		0.2		95	
14		70		1.5		48	
8		165		1.9		82	
10		90		2.1		50	
17		65		0.7		60	
20		60		1.1		55	
16		75		0.4		72	
12		110		1.2		60	
13		108		1.1		62	
11		112		1.3		59	
12		109		1.2		61	
12		111		1.1		60	
13		107		1.2		63	

# QXL DOE >Analyze Design >Run Regression



# 1 of 4 Regression Models in Coded Units

Inference Latency ms							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	Const	12.052	0.7514	16.039	0.000		
Batch Size Rows	Batch Size Rows (A)	-3.8	0.6887	-5.5177	0.000	1.0022	<input checked="" type="checkbox"/>
Quantization	Quantization (B)	4.887	0.6888	7.0952	0.000	1.0002	<input checked="" type="checkbox"/>
Precision Mode	Precision Mode (C)	1.5975	0.6888	2.3193	0.046	1.0002	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	AB	-0.9758	0.7696	-1.2679	0.237	1.0002	<input checked="" type="checkbox"/>
Batch Size Rows*Precision Mode	AC	-0.1898	0.7696	-0.2466	0.811	1.0002	<input checked="" type="checkbox"/>
Quantization*Precision Mode	BC	0.0	0.77	0.0	1.000	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization*Precision Mode	ABC	0.25	0.77	0.3247	0.753	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Batch Size Rows	AA	-1.2522	1.32	-0.9487	0.368	1.8204	<input checked="" type="checkbox"/>
Quantization*Quantization	BB	1.0	1.3133	0.7614	0.466	1.8182	<input checked="" type="checkbox"/>
Precision Mode*Precision Mode	CC	5.5	1.3133	4.1879	0.002	1.8182	<input checked="" type="checkbox"/>

## Hypothesis Test

**Ho: Coefficient is zero**

**Ha: Coefficient is not zero (term is significant)**

**Accept Ho if P-value > 0.10 (uncheck "In Model" box)**

## 2 of 4 Regression Models in Coded Units

Throughput inf/sec							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	Const	105.262	4.9538	21.249	0.000		
Batch Size Rows	Batch Size Rows (A)	38	4.5403	8.3695	0.000	1.0022	<input checked="" type="checkbox"/>
Quantization	Quantization (B)	-23.127	4.5408	-5.0932	0.001	1.0002	<input checked="" type="checkbox"/>
Precision Mode	Precision Mode (C)	-6.0446	4.5408	-1.3312	0.216	1.0002	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	AB	-9.5415	5.074	-1.8805	0.093	1.0002	<input checked="" type="checkbox"/>
Batch Size Rows*Precision Mode	AC	-3.347	5.074	-0.6596	0.526	1.0002	<input checked="" type="checkbox"/>
Quantization*Precision Mode	BC	5.625	5.0762	1.1081	0.297	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization*Precision Mode	ABC	1.875	5.0762	0.3694	0.720	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Batch Size Rows	AA	22.247	8.702	2.5566	0.031	1.8204	<input checked="" type="checkbox"/>
Quantization*Quantization	BB	-15.318	8.658	-1.7693	0.111	1.8182	<input checked="" type="checkbox"/>
Precision Mode*Precision Mode	CC	-25.318	8.658	-2.9243	0.017	1.8182	<input checked="" type="checkbox"/>

### Hypothesis Test

**Ho: Coefficient is zero**

**Ha: Coefficient is not zero (term is significant)**

**Accept Ho if P-value > 0.10 (uncheck "In Model" box)**

## 3 of 4 Regression Models in Coded Units

% Accuracy Drop							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	Const	1.2505	0.0762	16.421	0.000		
Batch Size Rows	Batch Size Rows (A)	0.14	0.0698	2.0059	0.076	1.0022	<input checked="" type="checkbox"/>
Quantization	Quantization (B)	-0.6403	0.0698	-9.1733	0.000	1.0002	<input checked="" type="checkbox"/>
Precision Mode	Precision Mode (C)	-0.7107	0.0698	-10.182	0.000	1.0002	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	AB	-0.024	0.078	-0.3074	0.766	1.0002	<input checked="" type="checkbox"/>
Batch Size Rows*Precision Mode	AC	-0.056	0.078	-0.7173	0.491	1.0002	<input checked="" type="checkbox"/>
Quantization*Precision Mode	BC	0.5	0.078	6.4075	0.000	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization*Precision Mode	ABC	0.0	0.078	0.0	1.000	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Batch Size Rows	AA	0.3605	0.1338	2.6946	0.025	1.8204	<input checked="" type="checkbox"/>
Quantization*Quantization	BB	0.0682	0.1331	0.5123	0.621	1.8182	<input checked="" type="checkbox"/>
Precision Mode*Precision Mode	CC	-0.5818	0.1331	-4.3715	0.002	1.8182	<input checked="" type="checkbox"/>

### Hypothesis Test

**Ho: Coefficient is zero**

**Ha: Coefficient is not zero (term is significant)**

**Accept Ho if P-value > 0.10 (uncheck "In Model" box)**

## 3 of 4 Regression Models in Coded Units

% CPU Utilization							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	Const	61.26	0.6591	92.948	0.000		
Batch Size Rows	Batch Size Rows (A)	15.8	0.6041	26.156	0.000	1.0022	<input checked="" type="checkbox"/>
Quantization	Quantization (B)	4.5866	0.6041	7.592	0.000	1.0002	<input checked="" type="checkbox"/>
Precision Mode	Precision Mode (C)	10.487	0.6041	17.359	0.000	1.0002	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	AB	-1.0058	0.6751	-1.4899	0.170	1.0002	<input checked="" type="checkbox"/>
Batch Size Rows*Precision Mode	AC	-0.9658	0.6751	-1.4307	0.186	1.0002	<input checked="" type="checkbox"/>
Quantization*Precision Mode	BC	1.0	0.6754	1.4807	0.173	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization*Precision Mode	ABC	-1.0	0.6754	-1.4807	0.173	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Batch Size Rows	AA	4.6491	1.1578	4.0157	0.003	1.8204	<input checked="" type="checkbox"/>
Quantization*Quantization	BB	-4.3182	1.1519	-3.7487	0.005	1.8182	<input checked="" type="checkbox"/>
Precision Mode*Precision Mode	CC	4.1818	1.1519	3.6304	0.005	1.8182	<input checked="" type="checkbox"/>

### Hypothesis Test

**Ho: Coefficient is zero**

**Ha: Coefficient is not zero (term is significant)**

**Accept Ho if P-value > 0.10 (uncheck "In Model" box)**

1) Uncheck the box for insignificant terms.

2) Reduce the model by running regression.

Inference Latency ms							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	Const	11.996	0.6134	19.557	0.000		<input checked="" type="checkbox"/>
Batch Size Rows	Batch Size Rows (A)	-3.8186	0.6121	-6.2386	0.000	1.0008	<input checked="" type="checkbox"/>
Quantization	Quantization (B)	4.9	0.6125	7.9997	0.000	1.0	<input checked="" type="checkbox"/>
Precision Mode	Precision Mode (C)	1.6	0.6125	2.6121	0.020	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	AB						<input type="checkbox"/>
Quantization*Precision Mode	BC						<input type="checkbox"/>
Batch Size Rows*Batch Size Rows	AA						<input type="checkbox"/>
Quantization*Quantization	BB						<input type="checkbox"/>
Precision Mode*Precision Mode	CC	5.3527	0.8666	6.1768	0.000	1.0008	<input checked="" type="checkbox"/>
	R <sup>2</sup>	0.9067					

1 of 4 Reduced Models in Coded Coefficients. Note the R-squared .

1) Uncheck the box for insignificant terms.

2) Reduce the model by running regression.

Throughput inf/sec							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	Const	103.373	5.0073	20.644	0.000		<input checked="" type="checkbox"/>
Batch Size Rows	Batch Size Rows (A)	38	4.6998	8.0854	0.000	1.0022	<input checked="" type="checkbox"/>
Quantization	Quantization (B)	-23.127	4.7004	-4.9203	0.000	1.0002	<input checked="" type="checkbox"/>
Precision Mode	Precision Mode (C)	-6.0	4.6998	-1.2766	0.224	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	AB	-9.5415	5.2523	-1.8167	0.092	1.0002	<input checked="" type="checkbox"/>
Quantization*Precision Mode	BC						<input type="checkbox"/>
Batch Size Rows*Batch Size Rows	AA	16.477	8.3513	1.973	0.070	1.5647	<input checked="" type="checkbox"/>
Quantization*Quantization	BB						<input type="checkbox"/>
Precision Mode*Precision Mode	CC	-31.063	8.3082	-3.7388	0.002	1.5625	<input checked="" type="checkbox"/>
	R <sup>2</sup>	0.8928					

2 of 4 Reduced Models in Coded Coefficients. Note the R-squared .

1) Uncheck the box for insignificant terms.

2) Reduce the model by running regression.

% Accuracy Drop							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	Const	1.2589	0.0648	19.427	0.000		<input type="checkbox"/>
Batch Size Rows	Batch Size Rows (A)	0.14	0.0608	2.3019	0.039	1.0022	<input checked="" type="checkbox"/>
Quantization	Quantization (B)	-0.64	0.0608	-10.523	0.000	1.0	<input checked="" type="checkbox"/>
Precision Mode	Precision Mode (C)	-0.71	0.0608	-11.674	0.000	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	AB						<input type="checkbox"/>
Quantization*Precision Mode	BC	0.5	0.068	7.3531	0.000	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Batch Size Rows	AA	0.3861	0.1081	3.5729	0.003	1.5647	<input checked="" type="checkbox"/>
Quantization*Quantization	BB						<input type="checkbox"/>
Precision Mode*Precision Mode	CC	-0.5563	0.1075	-5.1737	0.000	1.5625	<input checked="" type="checkbox"/>
	R <sup>2</sup>	0.9625					

3 of 4 Reduced Models in Coded Coefficients. Note the R-squared .

1) Uncheck the box for insignificant terms.

2) Reduce the model by running regression.

% CPU Utilization							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	Const	61.26	0.768	79.767	0.000		
Batch Size Rows	Batch Size Rows (A)	15.8	0.7039	22.447	0.000	1.0022	<input checked="" type="checkbox"/>
Quantization	Quantization (B)	4.6	0.7039	6.5352	0.000	1.0	<input checked="" type="checkbox"/>
Precision Mode	Precision Mode (C)	10.5	0.7039	14.917	0.000	1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	AB						<input type="checkbox"/>
Quantization*Precision Mode	BC						<input type="checkbox"/>
Batch Size Rows*Batch Size Rows	AA	4.6491	1.3491	3.4462	0.004	1.8204	<input checked="" type="checkbox"/>
Quantization*Quantization	BB	-4.3182	1.3423	-3.2171	0.007	1.8182	<input checked="" type="checkbox"/>
Precision Mode*Precision Mode	CC	4.1818	1.3423	3.1155	0.008	1.8182	<input checked="" type="checkbox"/>
	R <sup>2</sup>	0.9843					

4 of 4 Reduced Models in Coded Coefficients. Note the R-squared .

# QXL DOE >Analyze Design >Uncoded Coefficients

The screenshot displays the QXL DOE software interface. At the top, the title bar shows 'DOE 2 • Saved' and a search bar. The ribbon includes tabs for 'File', 'Home', 'Insert', 'Draw', 'Page Layout', 'Formulas', 'Data', 'Review', 'View', 'Automate', 'Help', 'QXL Monte Carlo', 'QXL Stat Tools', and 'QXL DOE'. The 'Analyze Design' group is expanded, showing 'Run Regression' and 'Uncoded Coefficients'. The 'Uncoded Coefficients' option is highlighted with a red box, and a tooltip is visible next to it. The tooltip text reads: 'Uncoded Coefficients Run regression in uncoded units from the design sheet or other regression sheet.'

AutoSave **On** DOE 2 • Saved

File Home Insert Draw Page Layout Formulas Data Review View Automate Help QXL Monte Carlo QXL Stat Tools **QXL DOE**

Design Create Modify Analyze Optimize Charts ProTest Options About  
Wizard Design Design Design Design (Covering Array) Options About

AA110

20  
21

**Uncoded Coefficients**  
Run regression in uncoded units from the design sheet or other regression sheet.

# 1 of 4 Regression Models in Uncoded Coefficients

Inference Latency ms							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	Const	16.324	0.9012	18.113			
Batch Size Rows	Batch Size Rows (A)	-0.5091	0.0816	-6.2386		1.0008	<input checked="" type="checkbox"/>
Quantization	Quantization (B)	4.9	0.6125	7.9997		1.0	<input checked="" type="checkbox"/>
Precision Mode	Precision Mode (C)	1.6	0.6125	2.6121		1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	AB						<input type="checkbox"/>
Quantization*Precision Mode	BC						<input type="checkbox"/>
Batch Size Rows*Batch Size Rows	AA						<input type="checkbox"/>
Quantization*Quantization	BB						<input type="checkbox"/>
Precision Mode*Precision Mode	CC	5.3527	0.8666	6.1768		1.0008	<input checked="" type="checkbox"/>
	R <sup>2</sup>	0.9067					

## 2 of 4 Regression Models in Uncoded Coefficients

Throughput inf/sec							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	Const	81.47	11.523	7.07			<input checked="" type="checkbox"/>
Batch Size Rows	Batch Size Rows (A)	0.0868	2.6234	0.0331		17.565	<input checked="" type="checkbox"/>
Quantization	Quantization (B)	-12.313	7.5295	-1.6354		2.5666	<input checked="" type="checkbox"/>
Precision Mode	Precision Mode (C)	-6.0	4.6998	-1.2766		1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	AB	-1.2722	0.7003	-1.8167		2.5666	<input checked="" type="checkbox"/>
Quantization*Precision Mode	BC						<input type="checkbox"/>
Batch Size Rows*Batch Size Rows	AA	0.2929	0.1485	1.973		18.298	<input checked="" type="checkbox"/>
Quantization*Quantization	BB						<input type="checkbox"/>
Precision Mode*Precision Mode	CC	-31.063	8.3082	-3.7388		1.5625	<input checked="" type="checkbox"/>
	R <sup>2</sup>	0.8928					

## 3 of 4 Regression Models in Uncoded Coefficients

% Accuracy Drop							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	<b>Const</b>	1.5962	0.1491	10.704			
Batch Size Rows	<b>Batch Size Rows (A)</b>	-0.098	0.0339	-2.8876		17.565	<input checked="" type="checkbox"/>
Quantization	<b>Quantization (B)</b>	-0.64	0.0608	-10.523		1.0	<input checked="" type="checkbox"/>
Precision Mode	<b>Precision Mode (C)</b>	-0.71	0.0608	-11.674		1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	<b>AB</b>						<input type="checkbox"/>
Quantization*Precision Mode	<b>BC</b>	0.5	0.068	7.3531		1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Batch Size Rows	<b>AA</b>	0.0069	0.0019	3.5729		18.298	<input checked="" type="checkbox"/>
Quantization*Quantization	<b>BB</b>						<input type="checkbox"/>
Precision Mode*Precision Mode	<b>CC</b>	-0.5563	0.1075	-5.1737		1.5625	<input checked="" type="checkbox"/>
	<b>R<sup>2</sup></b>	0.9625					

## 4 of 4 Regression Models in Uncoded Coefficients

% CPU Utilization							
Y-Hat							
Name	Factor	Coeff	SE	T	P	VIF	In Model
	<b>Const</b>	49.325	1.9086	25.844			
Batch Size Rows	<b>Batch Size Rows (A)</b>	0.7016	0.4216	1.6643		20.221	<input checked="" type="checkbox"/>
Quantization	<b>Quantization (B)</b>	4.6	0.7039	6.5352		1.0	<input checked="" type="checkbox"/>
Precision Mode	<b>Precision Mode (C)</b>	10.5	0.7039	14.917		1.0	<input checked="" type="checkbox"/>
Batch Size Rows*Quantization	<b>AB</b>						<input type="checkbox"/>
Quantization*Precision Mode	<b>BC</b>						<input type="checkbox"/>
Batch Size Rows*Batch Size Rows	<b>AA</b>	0.0827	0.024	3.4462		21.288	<input checked="" type="checkbox"/>
Quantization*Quantization	<b>BB</b>	-4.3182	1.3423	-3.2171		1.8182	<input checked="" type="checkbox"/>
Precision Mode*Precision Mode	<b>CC</b>	4.1818	1.3423	3.1155		1.8182	<input checked="" type="checkbox"/>
	<b>R<sup>2</sup></b>	0.9843					
	..						

# QXL Monte Carlo > Create/Modify Design Sheet > Import Model from Quantum XL DOE

AutoSave  On DOE 2 • Saved

File Home Insert Draw Page Layout Formulas Data Review View Automate Help **QXL Monte Carlo**

I/O Manager Mark Input Mark Output Unmark

**Design**

Create/Modify Design Sheet Move Cells Run Optimize Contribution Tools

Additional Tools

- Sampling Matrices Linear Tolerancing
- NOLHS Design
- Scorecard

- Custom Distributions
- Create Empirical
- Create Design Sheet
- Modify Design Sheet
- Merge Design Sheets
- Visual Merge
- Convert To Visual
- Import Model from Quantum XL DOE**
- Import Model from DOE Pro
- Validate Model
- Clean up workbook(s)

Parse DOE Regression Sheet ✕

Select outputs you want to import for Monte Carlo analysis.

Output	Import this output
Inference Latency ms Y-Hat	<input checked="" type="checkbox"/>
Throughput inf/sec Y-Hat	<input checked="" type="checkbox"/>
% Accuracy Drop Y-Hat	<input checked="" type="checkbox"/>
% CPU Utilization Y-Hat	<input checked="" type="checkbox"/>

Help Cancel Finish

# Enter 0 for 2<sup>nd</sup> Parameter in Monte Carlo Simulation Model

**Transfer function(s) in actual units.**

Process Inputs					Process Outputs				
Factor	Distro	1st Parameter	2nd Parameter	Exper	Name	Function	LSL	USL	Active
Batch Size Rows	Normal	8.5	0	8.5	Inference Latency ms	11.99634		11	X
Quantization	Normal	0	0	0	Throughput inf/sec	103.3726	125		X
Precision Mode	Normal	0	0	0	% Accuracy Drop	1.258867		2	X
Noise_Inference Latency ms	Normal	0	1.936972136	0	% CPU Utilization	61.25994		68	X
Noise_Throughput inf/sec	Normal	0	14.86222748	0	Use Quantum XL > Modify Design Sheet to add more outputs				
Noise_% Accuracy Drop	Normal	0	0.192328845	0					
Noise_% CPU Utilization	Normal	0	2.225880827	0					
Use Quantum XL > Modify Design Sheet to add more inputs									

**Firmware/Software parameters do not have tolerances.**

# QXL Monte Carlo > Optimize

Search

File Home Insert Draw Page Layout Formulas Data Review View Automate Help QXL Monte Carlo QXL Stat Tools QXL DOE

I/O Manager    Mark Input    Mark Output    Unmark    Create/Modify Design Sheet    Move Cells    Run    **Optimize**    Control

Design

Monte Carlo

Optimize

E8    X ✓    fx    0

### Optimizer

**Welcome**  
Welcome page

**Step 1 - Define Objectives**  
Define optimization objectives

**Step 2 - Define Constraints**  
Define optimization constraints

**Step 3 - Define Ranges**  
Define optimization ranges

**Step 4 - Define Weights**  
Define the dpm and weights

**Optimization Settings**  
Set the sample size

### Optimization

Use Optimization to reduce the defects per million of all the outputs simultaneously (Standard Mode) or to run an optimization on other statistics (e.g., mean, standard deviation, median, percentiles, etc.). Enter the objective, constraints, the range of valid inputs, and the type dpm to control the optimization.

**Step #1: Optimization Objective**  
Define the optimization objective(s). There can be more than one objective, but only one can be active at a time. If at least two outputs have a specification limit defined (upper or lower), run the optimization to minimize the dpm for all outputs.


**Step #2: Optimization Constraints**  
Optionally add optimization constraints to ensure optimization objectives are satisfied.

**Step #3: Define the Lows and Highs for the Input Variables**  
Define the low and high values for each input location parameter. The optimizer will find the optimal solution within the specified range.

**Step #4: Define the Type DPM and Weights**  
Optionally assign each output a type dpm and weight. The weight is only required if the model includes more than one output.

Help    < Previous    Next >    Save Settings    Cancel    Optimize

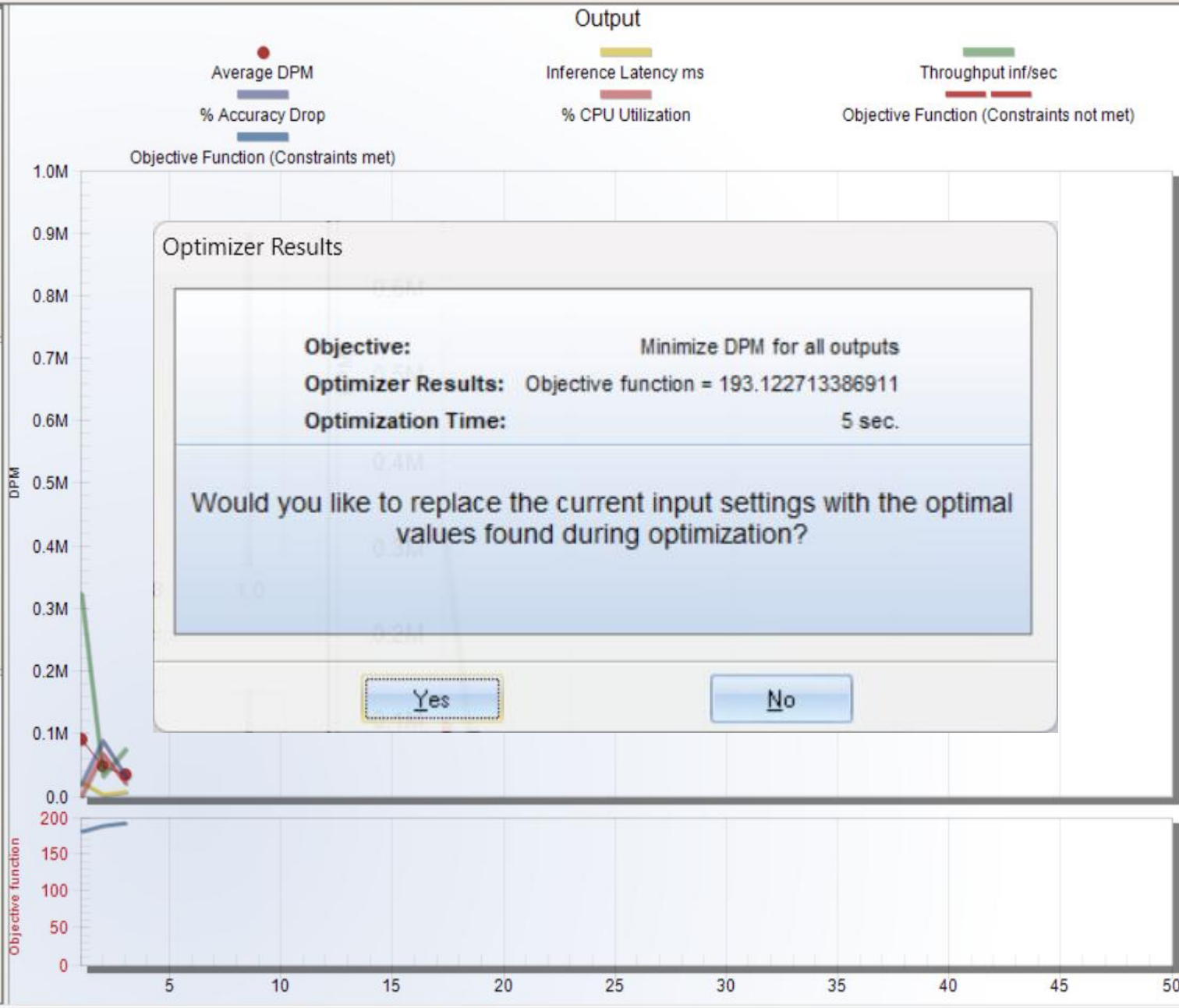
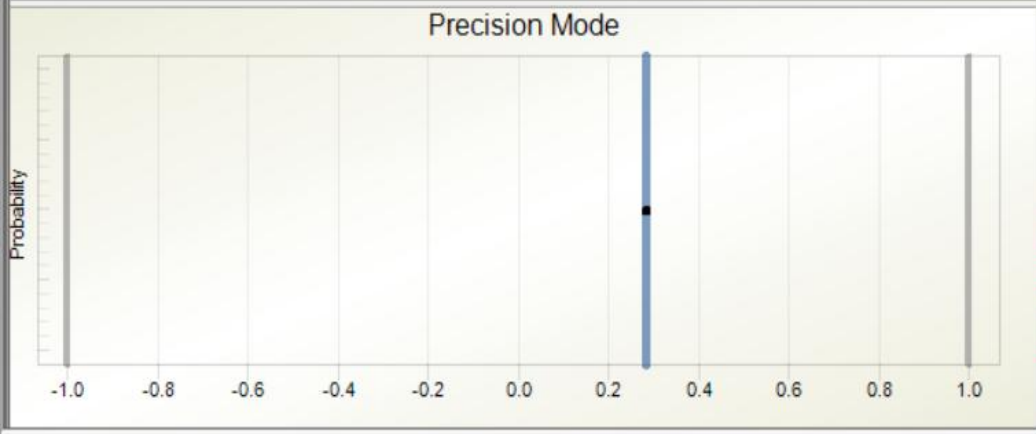
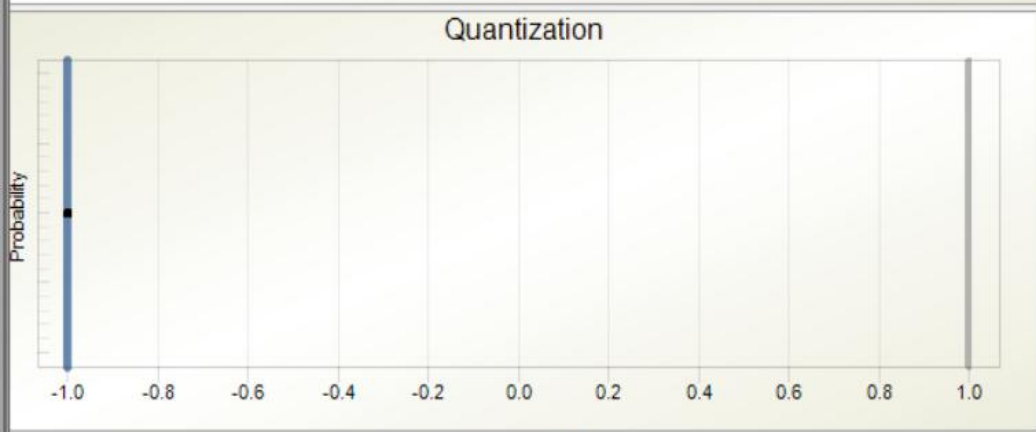
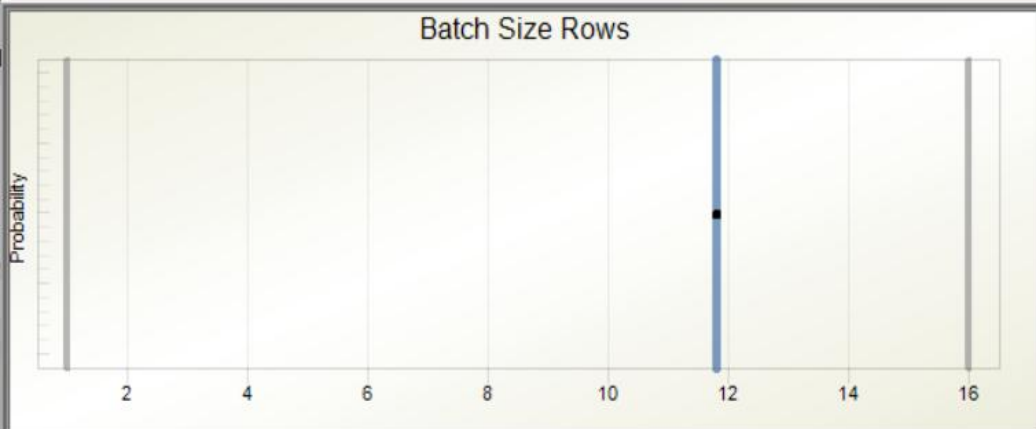
# Keep Default Settings for Optimizer Objectives

 **Click "Add Objective" to define the goal for the optimization.**  
Only one objective can be active at a time.

List of optimizer objectives	Active
<u>Minimize</u> the <u>DPM</u> of <u>all out...</u> <i>Minimize the DPM of all outputs</i>	<input checked="" type="checkbox"/>

Minimize  
Maximize  
Set target

All Outputs  
Inference Latency ms  
Throughput inf/sec  
% Accuracy Drop  
% CPU Utilization



# QXL Monte Carlo >Run >1,000,000 Simulations

The screenshot shows the QXL Monte Carlo software interface. At the top, there is a search bar and the text "Book4 - Excel". Below this is a ribbon with tabs: File, Home, Insert, Draw, Page Layout, Formulas, Data, Review, View, Automate, Help, and **QXL Monte Carlo**. The QXL Monte Carlo ribbon is divided into two sections: "Design" and "Additional Tools".

The "Design" section contains the following tools:

- I/O Manager
- Mark Input
- Mark Output
- Unmark
- Create/Modify Design Sheet
- Move Cells

The "Additional Tools" section contains the following tools:

- Run (selected, with a dropdown menu open)
- Optimize
- Contribution Tools
- Sampling Matrices
- Linear Tolerancing
- NOLHS Design
- Scorecard

The "Run" dropdown menu is open, showing the following options:

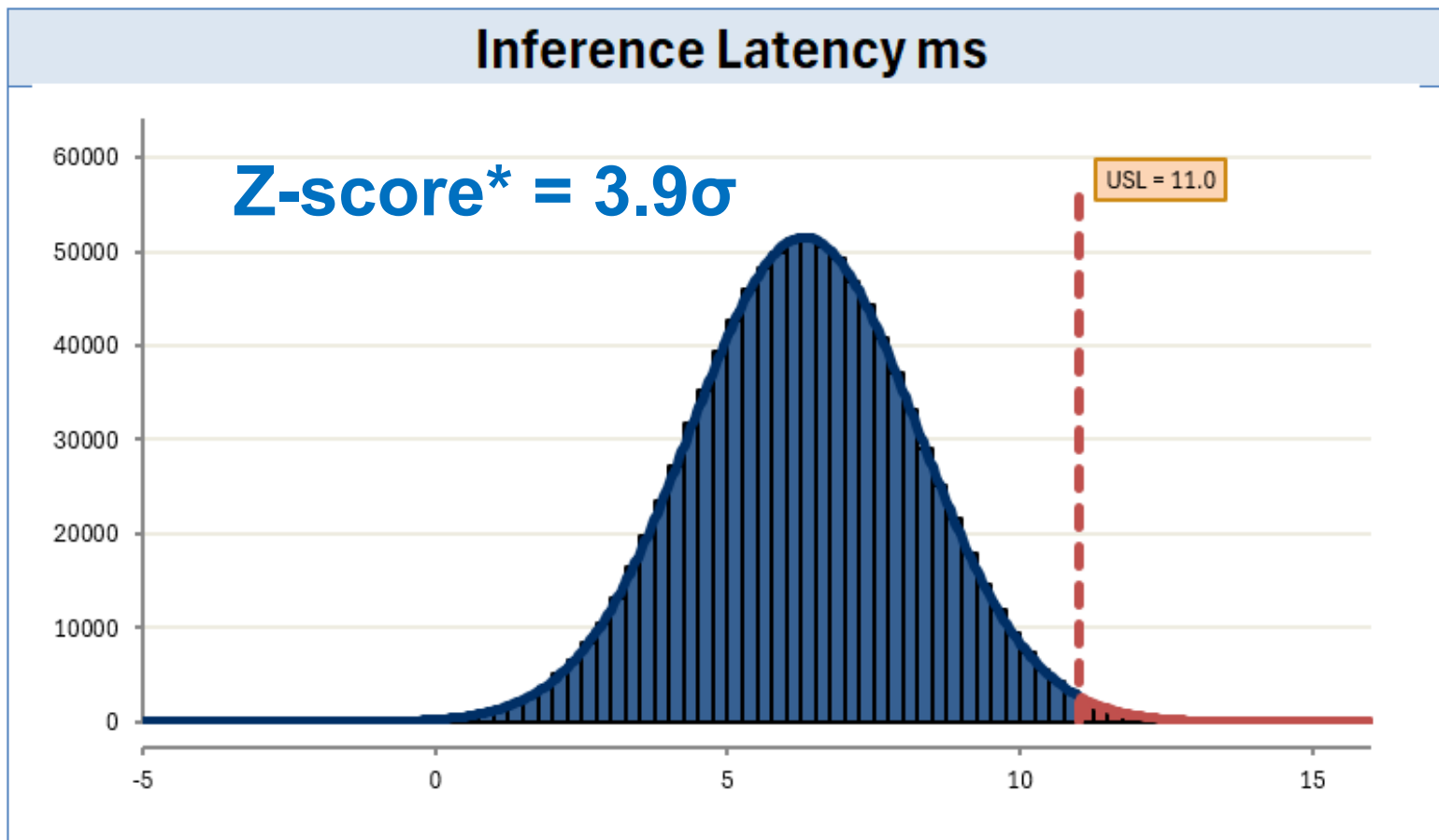
- 1,000 Simulations
- 10,000 Simulations
- 100,000 Simulations
- 1,000,000 Simulations** (highlighted with a mouse cursor)
- Custom**
- N Simulations...
- Set up Histogram Groups...

## Converting DPM to Z-score

- Defects Per Million (**DPM**) is a better indicator of quality than Cpk when:
  - The distribution is non-normal or asymmetrical.
  - One is using Monte Carlo simulation to generate a million data points.
- **Yield =  $(10^6 - \text{DPM}) / 10^6$**
- **Sigma Long-term =  $\text{NORM.S.INV}(\text{Yield})$**
- **Sigma Short-term or Z-score = Sigma Long-term + 1.5**

Process Inputs		
<b>Batch Size Rows</b>	Mean	StDev
Distribution: <b>Normal</b>	12	0
<b>Quantization</b>	Mean	StDev
Distribution: <b>Normal</b>	-1	0
<b>Precision Mode</b>	Mean	StDev
Distribution: <b>Normal</b>	0	0

Level	Quantization	Precision Mode
-1	INT4	FP16
0	INT8	TF32
+1	FP16	FP32



#### Process Outputs

# of Simulations	1,000,000
Mean	6.308709
StdDev	1.936854
Median	6.307316
LSL	
USL	11

#### Normal Distro Statistics

KS Test p-Value (Normal)	Not Avail
dpm	7,715
Cpk	0.807
Cp	

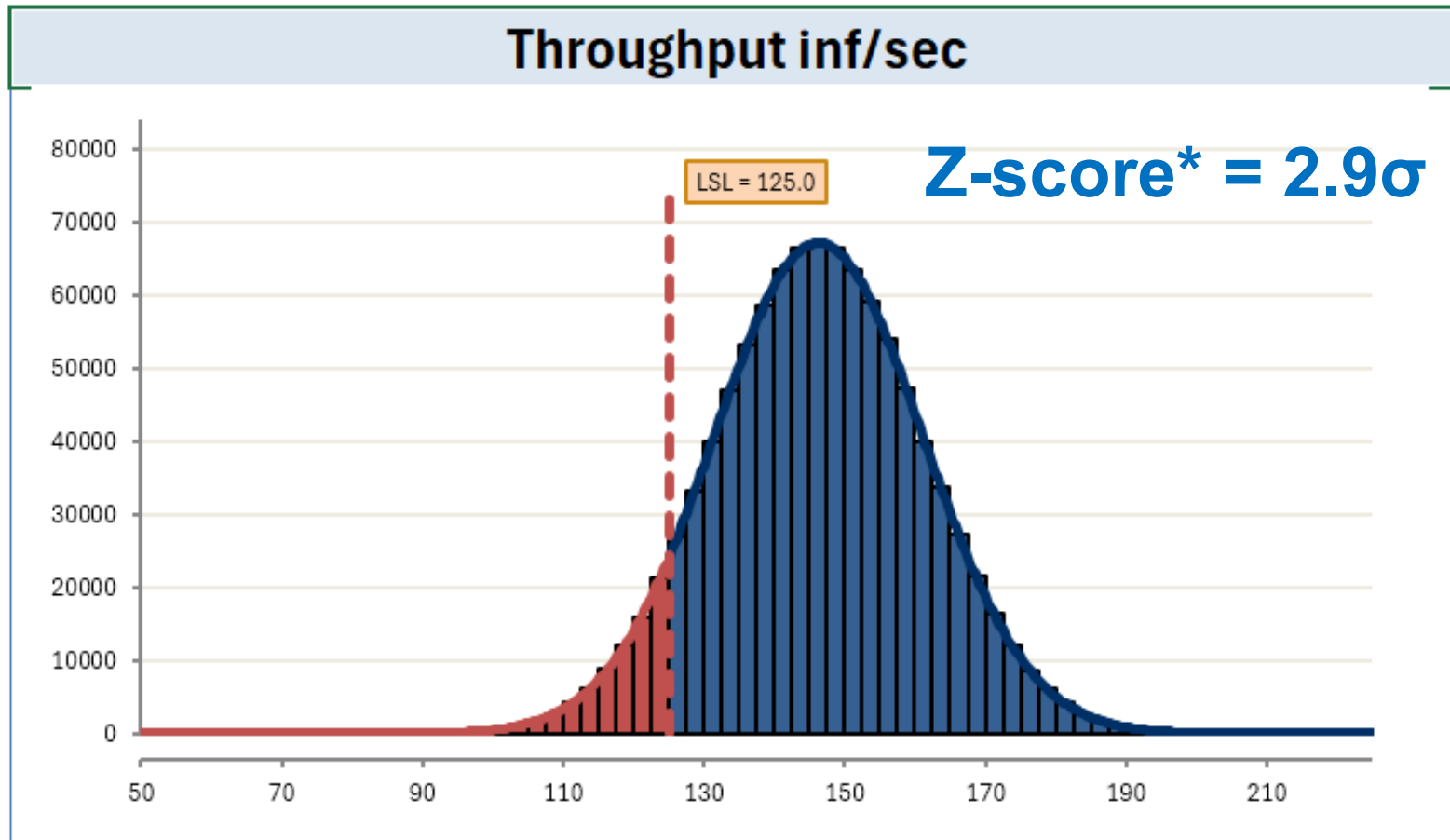
#### Observed Defect Statistics

Simulations outside of spec	7,795
-----------------------------	-------

\*note: based on Observed DPM

Process Inputs		
<b>Batch Size Rows</b>	Mean	StDev
Distribution: <b>Normal</b>	12	0
<b>Quantization</b>	Mean	StDev
Distribution: <b>Normal</b>	-1	0
<b>Precision Mode</b>	Mean	StDev
Distribution: <b>Normal</b>	0	0

Level	Quantization	Precision Mode
-1	INT4	FP16
0	INT8	TF32
+1	FP16	FP32



Process Outputs	
# of Simulations	1,000,000
Mean	146.29
StdDev	14.86
Median	146.30
LSL	125
USL	

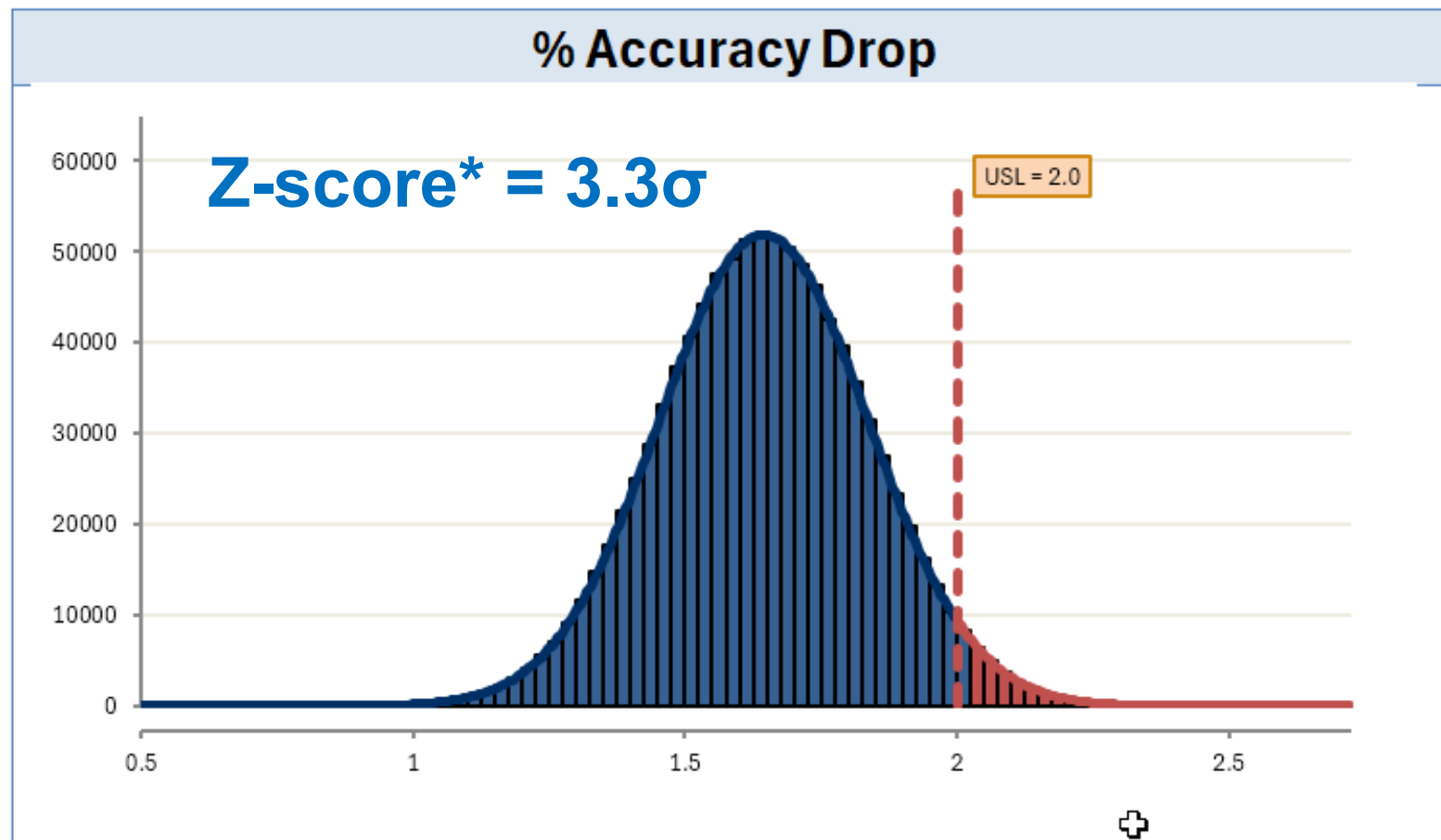
Normal Distro Statistics	
KS Test p-Value (Normal)	Not Avail
dpm	75,914
Cpk	0.478
Cp	

Observed Defect Statistics	
Simulations outside of spec	76,023

\*note: based on Observed DPM

Process Inputs		
<b>Batch Size Rows</b>	Mean	StDev
Distribution: <b>Normal</b>	12	0
<b>Quantization</b>	Mean	StDev
Distribution: <b>Normal</b>	-1	0
<b>Precision Mode</b>	Mean	StDev
Distribution: <b>Normal</b>	0	0

Level	Quantization	Precision Mode
-1	INT4	FP16
0	INT8	TF32
+1	FP16	FP32



Process Outputs	
# of Simulations	1,000,000
Mean	1.64
StdDev	0.19
Median	1.64
LSL	
USL	2

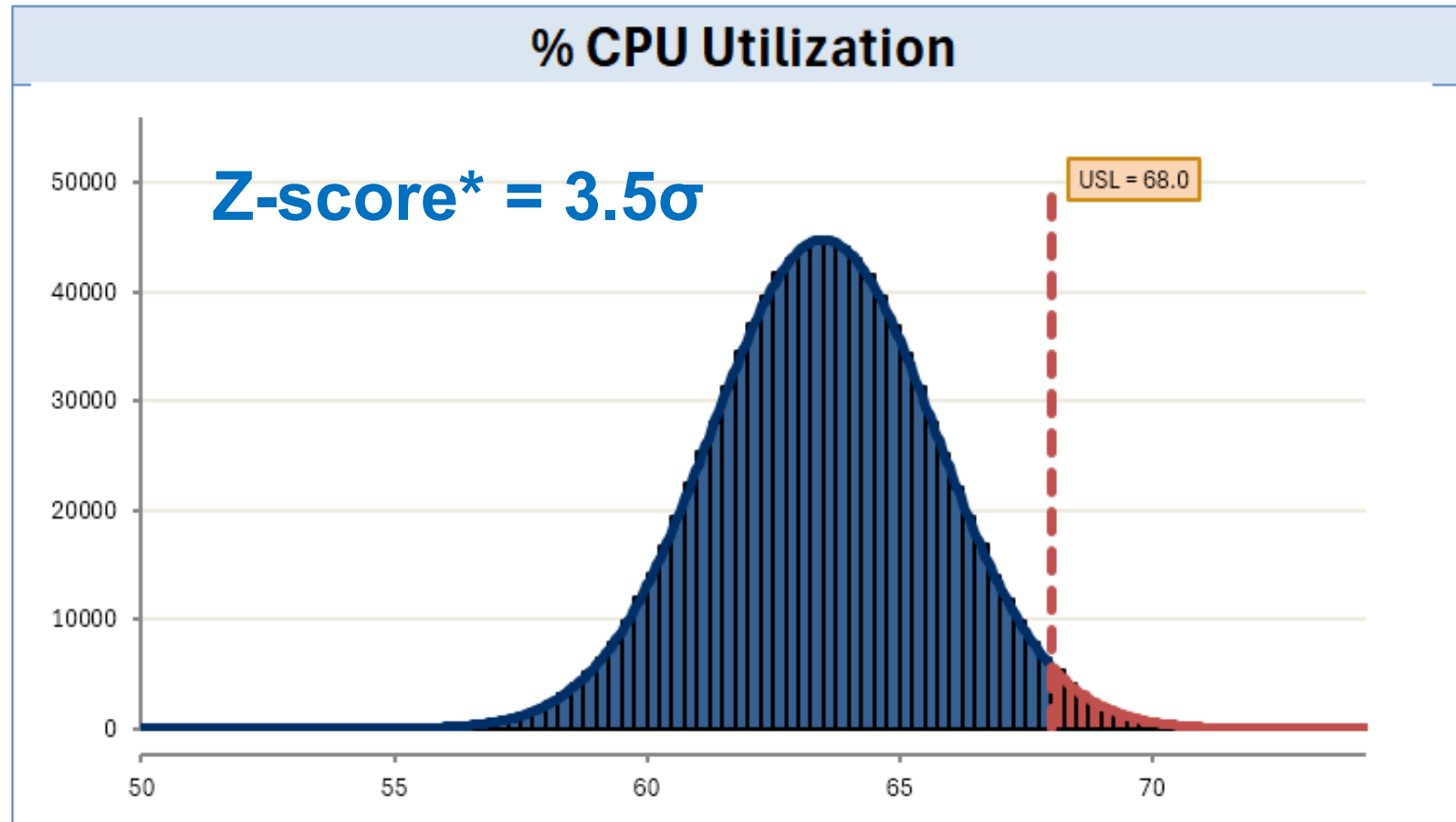
Normal Distro Statistics	
KS Test p-Value (Normal)	Not Avail
dpm	32,283
Cpk	0.616
Cp	

Observed Defect Statistics	
Simulations outside of spec	32,176

\*note: based on Observed DPM

Process Inputs		
<b>Batch Size Rows</b>	Mean	StDev
Distribution: <b>Normal</b>	12	0
<b>Quantization</b>	Mean	StDev
Distribution: <b>Normal</b>	-1	0
<b>Precision Mode</b>	Mean	StDev
Distribution: <b>Normal</b>	0	0

Level	Quantization	Precision Mode
-1	INT4	FP16
0	INT8	TF32
+1	FP16	FP32



Process Outputs	
# of Simulations	1,000,000
Mean	63.50
StdDev	2.23
Median	63.50
LSL	
USL	68

Normal Distro Statistics	
KS Test p-Value (Normal)	Not Avail
dpm	21,546
Cpk	0.674
Cp	

Observed Defect Statistics	
Simulations outside of spec	21,829

\*note: based on Observed DPM

# Summary

## OPTIMAL PARAMETERS

- **Batch Size = 12** per inference call
- **Quantization = INT4**
- **Precision Mode = TF32**

## CAPABILITY PER PERFORMANCE GOAL

- **Inference Latency (USL=11ms): 3.9 $\sigma$**
- **Throughput (LSL=125 ips): 2.9 $\sigma$**
- **Accuracy Drop (USL= 2%): 3.3 $\sigma$**
- **CPU Utilization (USL=68%): 3.5 $\sigma$**

## Conclusion:

- **Multi-objective optimization** was realized by creating four transfer functions based on empirical data.
- The Monte Carlo simulation served two purposes:
  1. It modeled the residuals, effectively **achieving an R-squared value of 100%**.
  2. It facilitated capability analysis using **one million samples**.
- **Quantum XL Version 18** is a powerful tool for **ML Inference Tuning**, as it integrates design of experiments with Monte Carlo simulation.