

■■■■■■■■ (Abstract)

AlphaClean (pipeline)

dataset

operators

(Quality Function)

■■■■ (Introduction)

Figure 1 illustrates the data flow and processing steps for the AlphaClean pipeline. The process begins with 'Data' and 'Scientist' inputs, which feed into a '60-80%' step. The output of this step is then processed by 'AlphaClean'. The final output is then handled by the 'operator' and 'pipeline' components.

■■■■■■■ (Methods)

AlphaClean ██████████: 1) ****Operators Library**** —
 remove duplicates, fill missing,
 normalize text 2) ****Pipeline Space Generator**** — ████████ pipeline
 operator ████████ 3) ****Quality Function**
(QF)** — ██████████ % missing < 1%, duplicates < 0.5% 4)
****Search Algorithm**** — ████ heuristic search ██ pipeline
 ██████████

■■■■■■■■■■ (Results)

AlphaClean dataset manual cleaning rule-based cleaning
: - clean - Pipeline
dataset

■■■■■■■ (Discussion)

****00000000:** - 00000000 000000000000 operator 000 - QF**
00000000000000 - Pipeline reuse 000 - 00 manual coding **00000000: -**

QF -
operator library - Dataset
:- dataset

 clean
 Master Schema

(Conclusion)

AlphaClean data cleaning pipeline

 pipeline reuse operator (OCR,
 ML/LLM