# Development of Python-based Toolkit to Improve Analysis of Atom Probe Tomography Data

Vishal Kanigicherla

March 2021

## Contents

**Abstract**

The purpose of the project is to develop a Python-based toolkit to optimize and streamline analysis of data produced by Atom Probe Tomography (APT), addressing many issues with the analysis process such as the difficulties of mass spectrum peak decomposition, mass spectrum ranging, and error calculation in reported values. This will be achieved using established Python libraries such as Pandas, processing large amounts of binary data generated by APT into workable CSV files, and addressing each issue in APT analysis with discrete scripts, finally compiling them into a toolkit available on Github.

# 1 Purpose

This project seeks to holistically address APT data analysis using Python, liberating APT data analysis from proprietary software while also improving existing data analysis techniques by implementing various algorithms employing optimization and machine learning. This project is trying to test the extent to which established techniques are user-friendly, accessible, and optimized.

# 2 Hypothesis

This project demonstrated the purpose by automating many menial tasks associated with APT data analysis, implementing methods to calculate important statistics integral to APT analysis, and tentatively will optimize industry standards for the specific algorithms involved in cluster identification and more. The automation was measured by consistency with past results of NRL APT data analysis, while the optimization will be measured against industry standards through efficiency and error comparison. The automation's consistency was a perfect equivalent to past calculations, and the spatial error in the opti-

mized algorithm will be statistically significantly less than that of the maximum separation algorithm.

# 3   Introduction

Atom probe tomography (APT) is a precise atom-by-atom dissection of a material volume, producing $250{\times}250{\times}1000$ nm$^3$ three-dimensional reconstructions, with sub-nanometer resolution, and chemical sensitivities approaching 10 atomic ppm (1). The ability to detect individual atomic species, directly in three dimensions, provides researchers with an unprecedented understanding of a material's nanostructure, which ultimately dictates that material's properties and performance (2).

Specimens for APT are in the form of sharply pointed needles, typically with an end radius less than 100 nm. These specimens are subjected to cryogenic temperature ( 30 K) and high voltage ( 10 kV), causing the atoms at the apex of the specimen tip to ionize and accelerate away from the positively charged tip and towards a detector. The detector records both the ion's time-of-flight (and hence its mass-to-charge ratio — its chemical identity) and its impact position (and hence its original location on the tip). APT specimens are generally prepared by milling in a dual-beam focused ion beam/scanning electron microscope (FIB/SEM) (2).

The ability to detect individual atomic species, directly in three dimensions, gives researchers an unprecedented understanding of a material's nanostructure, which ultimately dictates that material's properties and performance. APT also enables the intelligent design of future structural, electronic, optoelectronic, and functional materials for various applications.

# 4  Challenges of Atom Probe Tomography

Though APT is an incredibly powerful analytical technique, often considered a "holy grail" of such techniques, it has its downsides. FIB milling implants gallium ions into the specimen which can introduce damage and must be corrected for during analysis (4). APT also has a relatively small field of view, limiting the size of the nanoscale features that can be analyzed (3). Also, there are many artifacts associated with the technique — including local magnification and isotopic overlaps in the mass spectrum — which can complicate data analysis and interpretation. Developing new algorithms to overcome some of these challenges is one of the objectives of this proposed project.

Data analysis automation and optimization is another challenge with APT data, especially considering that modern APT datasets can contain hundreds of millions of atoms. One process which would benefit from automation is mass spectrum peak identification, which becomes cumbersome due to ambiguities in mass-to-charge ratio. A common example in steels occurs at 27 Da, which could be 27Al1+, 54Fe2+, or 54Cr2+. This ambiguity of isotopic overlap is normally resolved through user-interaction, but could easily be automated considering the natural abundances of the atomic species in question and their non-overlapping peaks. The primary data format of APT data is the POS file, which contains position (x, y, z) and mass-to-charge ratio for each atomic data point, and an accompanying range file, which identifies and assigns element information to peaks on a mass spectrum. New tools will be developed to work with these data files in Python. A common analytical technique in APT is the proximity histogram, or "proxigram," which plots the concentrations of elements with respect to a three-dimensional isoconcentration surface in the data. The proxigram is commonly used in processes of precipitation, where the precipitates can be delineated by an isoconcentration surface and then co-positions of the precipitates and the surrounding matrix can be computed. All of these tasks can be streamlined into one toolkit, again simplifying many menial processes

central to APT analysis.

There are other issues that may be addressed with this toolkit. These include implementing algorithms for detecting solute clusters, which are based on the maximum separation algorithm. This process, however, is heavily dependent on user-identified parameters such as dmax (maximum distance between atoms in a cluster), Nmin (minimum number of atoms in a cluster), and K, the order of the nearest-neighbour model sometimes used to determine dmax (5). Recent research has identified a possible avenue of optimization of this rather sensitive and subjective process utilizing Ripley's K-function and machine learning, but the problem remains open for further investigation (6). Other problems with APT data analysis include the optimization of existing simulations for field evaporation behavior such as TAPSim; while fruit for further discussion, may lay outside the scope of the current project (7).

The application of APT is extremely important in materials science. Research using the technique extends to diverse areas of interest, such as high-entropy alloys that have application as substitutions for titanium in corrosive environments, the characterization of biological metal-interacting compounds such as porphyrin rings or proteins such as ferritin, and geosciences and investigation of the makeup of our planet (8, 9, 1). The significance of this work is that implementation of such a toolkit will allow for an easily accessible and optimized data analysis procedure with a variety of scripts and algorithms necessary for APT data analysis, giving researchers the flexibility to turn their attention more towards the impact of their research itself, rather than the menial tasks associated with the tool that is APT. In stating such, the proposed toolkit will improve efficiency in the fields of research associated with APT, and will impact the diverse and extensive fields of research that APT finds application in.

# 5   Intended Outcome

The intended outcome of this work will be measured in two ways. First, for the sections of the toolkit intended to simply streamline certain data analysis processes, said sections will process past data collected by the NRL, and resultant statistics should be found consistent with past reports. For the exploratory sections of the toolkit with no prior implementation in NRL work, I will compare statistical distributions of data processed using new algorithms with distributions I generate through the implementation of standard algorithms using necessary two-sample T-tests and one-way ANOVAs as necessary.

# 6   Materials

Materials that were used for this project include Python libraries (pandas, numpy, matplotlib, tkinter, struct) to interact with files and data types for processing, a laptop with Python and necessary environment installed (currently Lenovo Thinkpad) including IVAS LT, a text editor (VSCode), and an HDF file viewer. Finally, I made use of POS/EPOS files, range files, and proxigram files publicly available or generated from past and current work at the NRL, in addition to Excel workbooks with formulae kindly provided by Dr. Colin Stewart at the NRL MSTD.

# 7   Procedure

Python code was written in order to:

- Generate initial CSV files from proxigram XLSX.

- Take user-input for manual peak decomposition in proxigrams.

- Convert proxigram atom counts to at% values.

- Generate profiles of proxigram data after discarding irrelevant isotopes.

- Create core statistics file from user identification of matrix and precipitate of profile.

- Turn large POS files into workable CSVs.

- Calculate and display spatial error statistics for proxigram and CSV.

- Take inputs for calculating average radius, vol fraction, and number density, do the calculation, and report the uncertainty.

- Generate a mass spectrum graph and analysis from POS CSV.

- Consolidate programs into GUI using Py2App and tkinter for better user experience.

# 8 CSV Generation and Manipulation

The *pandas* library is a powerful tool that can be used to manipulate large amounts of data in comma-separated values, or CSV files. One other industry standard file format is Hierarchical Data Format (HDF4, HDF5), though interaction with and export capability is also integrated into pandas.

In order to streamline the Excel data manipulation, proxigram files were exported as CSVs. POS and EPOS files, though, do not have this built in capability. Accordingly, a program was written to parse the binary formatting of reference POS files using struct, and write the 16 or 44 char string (.pos and .epos respectively) to a CSV, to be processed by later programs.

A bulk of the code was written for interaction with proxigrams, due to their easier manageability as a result of smaller file sizes. One concern, as reflected earlier, was the decomposition of overlapping peaks, previously being

resolved at the NRL with a manual input of known isotope abundances and requiring much work to conduct. The script written to automate much of this process interacts with three files. The first is the original CSV proxigram, from which unknown peaks are separated into their constituent elements, and the file is rewritten into an intermediate. This intermediate file is processed and resolved but still a standard proxigram file, making it extremely relevant for other calculations. The third file is the final, with percentages instead of atom counts, and data deemed irrelevant to analysis by the user discarded; that file is used to generate the scatter plot to visualize concentration distribution, while the intermediate second file will be used after 'core' region identification for further statistics. Scripts automating many of these subtasks which also included a proxigram profile plotter, core composition statistics generation, spatial error CSV generation, and more were also compiled into a master script in the repository as *ProxigramPeakDecomp.py* and also divided into constituent tasks and placed in the repository.

## 9 Gibbsian Interfacial Excess of Solute

The Gibbsian interfacial excess of solute ($\Gamma_s$) is a quantity that relates the concentration of solute within an infinitesimally small interface, a "Gibbs dividing surface," between two bulk phases of a material. The data returned by standard atom probe tomography (position and mass-to-charge ratio), because it directly displays atom counts, can be easily used for calculating $\Gamma_s$ with respect to an arbitrary interface that is curved.

In many prior formulations, knowledge of the area of the interface is necessary to computing $\Gamma_s$. This often results in inherently error-producing somewhat arbitrary estimations of curves and planes. This is resolved by calculating area through other statistics that can be generated from APT: shell thickness, atomic density, and number of atom counts in each slice. Hellman

& Seidman (2002) provides a formula using summations to calculate $\Gamma_s$ from concentrations in shells from the interface of interest (10).

The purpose of a proximity histogram, or proxigram (ubiquitous across APT research) is to visualize atom counts in the vicinity of an isoconcentration surface, usually a nonplanar interface of interest. The proxigram, therefore, is ideal for calculating $\Gamma_s$ from the methods established in this article. Summations ($\sum$) are easily executable in code through for-loops. Thus, I wrote a script to parse through a proxigram and sum concentrations of solutes of interest, and use statistics calculated in other package scripts (number density, particle volume, percent concentrations, etc.) to implement the formulae devised. The advantage of the formulation presented by Hellman & Seidman (2002) was that it eliminated the need for *a priori* knowledge of slice area, and through dimensional analysis substituted a formula for it (10).

$$A_i = \frac{N_i}{\rho \Delta l}$$

The difficulty arises when calculating an atomic number density $\rho$. This was circumvented by prompting the user for the lattice parameter of the element of interest, and assuming a Bravais lattice geometry (possibly open for input), we calculated the ideal atomic number density and accordingly approximated the area of the unit of interest. Ultimately, I was able to develop a script to calculate the important thermodynamic quantity of Gibbs interfacial excess of solute. The script for this is available in the Github repository as *GibbsInterfacialExcess.py*.

## 10    Results and Discussion

All relevant documentation along with sample data sets, errors, and inputs and outputs are available on a Github repository at the link https://github.com/sakanak/apt-

csv-work. The README will soon be updated to reflect more details about specific aspects of each file, the prerequisite information required, as well as helpful tips for users. At the moment, all algorithms and routines written have been mainly for the purpose of automation, streamlining the process of menial tasks: because the programs written complete their tasks of automation, reorganization, and calculation the results are that all goals have been met to their intended degree.

In a discussion of the implemented formulation of Gibbs interfacial excess of solute, we must address the use of prompting lattice parameters and the underlying assumptions that follow. The $\rho$ term in our use is for an extremely specific case; the spatial error may be calculated and revealed to negligible, yet it is still relevant to address the assumption of Bravais lattice structure in number density ($\rho = \frac{lattice\ points}{lattice\ parameter^3}$). We can begin to address the specifics by referencing the initial formulations for Krakauer & Seidman et al. (1993) for an APFIM TEM calculation of Gibbsian excess as well as Cahn's formulation, which does have limited application, to achieve a more rigorous description of $\rho$ to calculate the excess (10, 11).

The question of user experience is also an important one, which has been addressed in brief but merits further consideration. The best choice is always for there to be complete automation, no work or *a priori* considerations necessary for use of programs (including experience in chemical analysis or computer science). While this could not be circumvented in earlier stages of the project, it was instead addressed by clear, thorough prompts for value inputs. The current interface application works by interacting with the terminal. In conjunction with the README, while there is much improvement to be made, this will function as a stopgap until consolidation, whether into a package with documentation or a multi-platform compliant applet, is completed.

# 11    Conclusion

After developing the toolkit outlined in the proposal, we wanted to answer the question of whether APT data analysis can be improved with respect to modularity, ease-of-use, streamlining, and optimization. The value of this project is in the consolidation and improvement of established analytical tools as well as the development of new analysis techniques, resulting in improvements to any field of research with the viability of using such a technique. Tentative extensions to this toolkit include the implementation of alternatives to the standard maximum separation algorithm used in cluster identification and an exploration of atom probe tip behavior simulations.

# 12    In Progress

There are many tasks for which automation is a possibility in the realm of Python-based APT data analysis. First, the question of consolidation and packaging of scripts must be answered. Currently, I am developing both a Python package as well as an independent applet through which users may work with proxigram and POS/EPOS CSV data files. The scripts that I am working on at the moment and will tentatively complete by the end of June 2021 are listed below. This is a complete list of next steps in and of itself, but higher priorities moving forward.

- Optimizing cluster detection

- Generating volume and radii calculations with particle info

- Executing peak detection and automatic ranging and isotope assignment using available isotope abundances table

- Modelling matrix composition to track precipitation as a function of time

- *Additional tasks are to be determined*

# 13   Author's Note

Due to the nature of the current COVID-19 pandemic, it is difficult to develop research consistent with past APT research conducted through the Mentorship program at the NRL. Thank you for your continued patience and understanding during this time.

# 14   References