



---

Development and Semi-Automated Analysis of an *in vitro* Dissemination Model  
for Myeloma Cells Interacting with Mesenchymal Stromal Cells

Entwicklung und semi-automatisierte Analyse eines *in vitro* Modells  
für die Disseminierung von Myelomzellen in Interaktion mit mesenchymalen Stromazellen

---

Doctoral Thesis for a Doctoral Degree

*at the*

GRADUATE SCHOOL OF LIFE SCIENCES,  
JULIUS-MAXIMILIANS-UNIVERSITÄT WÜRZBURG,  
SECTION BIOMEDICINE

*submitted by*

**Martin Kuric**

*from*

Bad Neustadt a.d. Saale

Würzburg, 2024



**Submitted on:** .....  
Office stamp

**Members of the *Promotionskomitee*:**

**Chairperson:** Prof. Dr. Uwe Gbureck  
**Primary Supervisor:** Prof. Dr. rer. nat. Regina Ebert  
**Supervisor (Second):** Prof. Dr. med. Franziska Jundt  
**Supervisor (Third):** Prof. Dr. rer. nat. Torsten Blunk

**Date of Public Defense:** .....

**Date of Receipt of Certificates:** .....

This work was conducted at the Department of Musculoskeletal Tissue Regeneration (Bernhard-Heine-Centre for Locomotive Research), University of Würzburg from 08.04.2018 to 31.03.2024 under the supervision of Prof Dr. rer. nat. Regina Ebert.

---

## Acknowledgements

Lorem Ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum

## **Summary**

  Lorem Ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum

## **Zusammenfassung**

  Lorem Ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum bla

# Contents

<b>Summary</b>	<b>ii</b>
<b>Introduction</b>	<b>1</b>
Human Mesenchymal Stem/Stromal Cells . . . . .	2
Multiple Myeloma . . . . .	3
Myeloma-hMSC Interactions . . . . .	3
Myeloma Bone Disease . . . . .	3
Dissemination of Myeloma Cells . . . . .	4
Code-Automation as a Standard in Modern Biosciences . . . . .	5
How Code Quality Improves Scientific Reproducibility . . . . .	6
Python as a Programming Language . . . . .	8
Data Science with Python . . . . .	12
Aims . . . . .	14
<b>Chapter 1: Modelling Myeloma Dissemination <i>in vitro</i></b>	<b>15</b>
Abstract . . . . .	15
Introduction . . . . .	15
Methods . . . . .	15
Results . . . . .	15
Discussion . . . . .	15
Supplementary Figures and Methods . . . . .	16
<b>Chapter 2: Semi-Automation of Data Analysis</b>	<b>52</b>
Abstract . . . . .	52
Introduction . . . . .	52
Statement of Need . . . . .	55
Example . . . . .	55
Overview . . . . .	57
Discussion . . . . .	58
<b>Summarising Discussion</b>	<b>61</b>
Time Lapse . . . . .	61
Myeloma . . . . .	61
Semi-Automated Analysis Improves Agility During Establishing new <i>in vitro</i> Methods . . . . .	61
<b>References</b>	<b>63</b>
<b>Appendix A</b>	<b>68</b>
Class Diagram of plotastic . . . . .	69
Readme of plotastic . . . . .	72
Example Analysis “qpcr” using plotastic . . . . .	85
<b>Appendix B</b>	<b>91</b>
Author Contributions to Research Projects . . . . .	92
Author Contributions to Figures and Tables . . . . .	95
Affidavit . . . . .	98
Curriculum Vitae . . . . .	99

## **Introduction**

To provide a comprehensive background for the following chapters that focus on the interaction of human mesenchymal stromal cells (hMSCs) with multiple myeloma (MM) cells, this

## Human Mesenchymal Stem/Stromal Cells

Explaining what a mesenchymal stromal cell (MSC) is, is not such an easy task as one might expect. MSCs are derived from multiple MSCs different sources, serve a wide array of functions and are always isolated as a heterogenous group of cells. This makes it particularly challenging to find a consensus on their exact definition, nomenclature, exact function and *in vivo* differentiation potential. Therefore, the most effective approach to describe hMSCs is to present their historical context.

hMSCs first gained popularity as a stem cell. Stem cells lay the foundation of multicellular organisms. Embryonic stem cells orchestrate the growth and patterning during embryonic development, while adult stem cells are responsible for regeneration during adulthood. The classical definition of a stem cell is that of a relatively undifferentiated cell that divides asymmetrically, producing another stem cell and a differentiated cell (Cooper, 2000; Shenghui et al., 2009). Because of their significance in biology and regenerative medicine, stem cells have become a prominent subject in modern research. Especially human mesenchymal stromal cells (hMSCs) have proven to be a promising candidate in this context (Ullah et al., 2015).

Mesenchyme first appears in embryonic development during gastrulation. There, cells that are committed to a mesodermal fate, lose their cell junctions and exit the epithelial layer in order to migrate freely. This process is called epithelial-mesenchymal transition (Tam & Beddington, 1987; Nowotschin & Hadjantonakis, 2010). Hence, the term mesenchyme describes non-epithelial embryonic tissue differentiating into mesodermal lineages such as bone, muscles and blood. Interestingly, it was shown nearly twenty years earlier that cells within adult bone marrow seemed to have mesenchymal properties as they were able to differentiate into bone tissue (A. J. Friedenstein et al., 1966; A. Friedenstein & Kuralesova, 1971; Bianco, 2014). This was the origin of the “mesengenic process”-hypothesis: This concept states that mesenchymal stem cells serve as progenitors for multiple mesodermal tissues (bone, cartilage, muscle, marrow stroma, tendon, fat, dermis and connective tissue) during both adulthood and embryonic development (A. Caplan, 1991; A. I. Caplan, 1994). The mesenchymal nature of these cells (termed bone marrow stromal cells: BMSCs) was confirmed later when they were shown to differentiate into adipocytic (fat) and chondrocytic (cartilage) lineages (Pittenger et al., 1999). Since then, the term “mesenchymal stem cell” (MSC) has grown popular as an adult multipotent precursor to a couple of mesodermal tissues. hMSCs derived from bone marrow (hMSCs) were shown to differentiate into osteocytes, chondrocytes, adipocytes and cardiomyocytes (Gronthos et al., 1994; Muruganandan et al., 2009; Xu et al., 2004) Most impressively, these cells also exhibited ectodermal and endodermal differentiation potential, as they produced neuronal cells, pancreatic cells and hepatocytes (Barzilay et al., 2009; Wilkins et al., 2009; Gabr et al., 2013; Stock et al., 2014).

Furthermore, cultures with MSC properties can be established from “virtually every post-natal organs and tissues”, and not just bone marrow (da Silva Meirelles et al., 2006). However, it has to be noted that hMSCs can differ greatly in their transcription profile and *in vivo* differentiation potential depending on which tissue they originated from (Jansen et al., 2010; Sacchetti et al., 2016).

Since “hMSCs” are a heterogenous group of cells, they were defined by their *in vitro* characteristics. A minimal set of criteria are the following (Dominici et al., 2006): First, hMSCs must be plastic adherent. Second, they must express or lack a set of specific surface antigens (positive for CD73, CD90, CD105; negative for CD45, CD34, CD11b, CD19). Third, hMSCs must differentiate to osteoblasts, adipocytes and chondroblasts *in vitro*. Together, hMSCs exhibit diverse differentiation potentials and can be isolated from multiple sources of the body. This offers great opportunity for regenerative medicine, if the particular hMSC-subtype is properly characterized.

## Multiple Myeloma

Multiple myeloma arises from clonal expansion of malignant plasma cells in the bone marrow (BM). At diagnosis, myeloma cells have disseminated to multiple sites in the skeleton and, in some cases, to virtually any tissue (Rajkumar & Kumar, 2020; Bladé et al., 2022).

## Myeloma-hMSC Interactions

Since plasma cells can not survive outside the bone marrow, MM cells also require survival signals for growth and disease progression. These signals are produced by the bone marrow microenvironment, including ECM, MSCs and ACs (Kibler et al., 1998; García-Ortiz et al., 2021).

## Myeloma Bone Disease

Bone is a two-phase system in which the mineral phase provides the stiffness and the collagen fibers provide the ductility and ability to absorb energy (Viguet-Carrin et al., 2006). On a molecular level, bone tissue is composed of extracellular matrix (ECM) proteins that are calcified by hydroxyapatite crystals. This ECM consists mostly of collagen type I, but also components with major regulatory activity, such as fibronectin and proteoglycans that are essential for healthy bone physiology (Alcorta-Sevillano et al., 2020). Bone tissue is actively remodeled by bone-forming osteoblasts and bone-degrading osteoclasts. Osteoblasts are derived from mes-

enchymal stromal cells (MSCs) that reside in the bone marrow (A. J. Friedenstein et al., 1966; Pittenger et al., 1999). MSCs also give rise to adipocytes (ACs) to form Bone Marrow Adipose Tissue (BMAT), which can account for up to 70% of bone marrow volume (Fazeli et al., 2013).

MM indirectly degrades bone tissue by stimulating osteoclasts and inhibiting osteoblast differentiation, which leads to MM-related bone disease (MBD) (Glavey et al., 2017). MBD is present in 80% of patients at diagnosis and is characterized by osteolytic lesions, osteopenia and pathological fractures (Terpos et al., 2018).

## **Dissemination of Myeloma Cells**

dissemination is still widely unclear - multistep process - invasion, intravasation, intravascular arrest, extravasation, colonization - overcome adhesion, retention, and dependency on the BM microenvironment - loss of adhesion factors such as CD138

## Code-Automation as a Standard in Modern Biosciences

Beschreibe die Situation. - Big Data in Biosciences - what is big data, examples - Define citable challenges: - reproducibility crisis - lack of tools

In recent years, the biosciences have evolved dramatically, with a marked increase in the volume and complexity of data generated (Yang et al., 2017; Ekmekci et al., 2016). This transformation necessitates robust software tools, many of which require coding skills to use effectively. Here we summarize standard tools used by biosciences today and show their reliance on coding. The author argues that the role of a modern independent researcher is now intertwined with coding skills similar to a role of “precision medicine bioninformatician” (Gómez-López et al., 2019).

Statistical analysis in biosciences has traditionally been reliant on user-friendly tools like Excel and GraphPad Prism. While Excel by itself is recognized as limited for complex data analysis (Tanavalee et al., 2016; Incerti et al., 2019), GraphPad Prism offers more advanced statistical models .

However, increasingly demands more sophisticated approaches as data sets grow in size and complexity.

R and Python scripts offer more efficient and versatile solutions, enabling complex analyses with a few lines of code (R Core Team, 2018; Vallat, 2018).

Recognizing this trend, Microsoft has integrated a Python interpreter into Excel to computations more accessible within a widely used platform (?).

Next-generation sequencing, such as bulk RNAseq, has become affordable, allowing for larger sample sets during a single PhD project. This technology offers advanced tools that are most efficiently used through scripting in R or Python. In the absence of a dedicated statistician, researchers are compelled to learn coding.

In gene ontology, tools such as Metascape facilitate the integration of vast datasets and outputs multiple useful data visualizations. Metascape also provides multiple excel sheets, containing all results, sometimes in a nested format, which provides even further information that's adaptable for specific hypotheses, but given the sheer amount of data, is impractical to analyze manually.

since Metascape returns large Excel sheets with complex nested information, a researcher without coding skills requires manual work to adapt the results to specific research hypotheses.

its true potential is unlocked only when researchers can manipulate and analyze these data through scripting.

Modern gene ontology tools like Metascape offer powerful graphical user interfaces. However, their effectiveness is only possible through standardizing multiple large datasets.

The output from Metascape, large Excel sheets with complex nested information, is more efficiently analyzed through scripting, which is often necessary to adapt metascape results to specific research hypotheses.

Image analysis is another area where coding skills are essential. ImageJ/FIJI, a standard tool in the field, requires scripting for batch processing of multiple images and automating multiple processing steps into a pipeline. While macros can be recorded, understanding the underlying code is necessary for troubleshooting and adapting the macro to new datasets.

In the field of protein structural biology, Pymol is a standard tool that also has a Python command interface.

Similarly, artificial intelligence (AI), a game-changer in biomedicine, primarily uses Python due to its extensive libraries for machine learning and scientific computing. Python is also a standard for integrative biomedicine simulations.

Finally, databases and repositories are essential for storing, retrieving, and sharing data. Researchers need to understand common file formats to adhere to standards that ensure re-usability and interoperability. Scripting helps automate the process of formatting data for submission to these databases.

In conclusion, the integration of coding in bioscience research is not just a trend but a necessity. As the field continues to evolve, the demarcation between biologists and computational scientists blurs, underscoring the importance of coding skills for the next generation of researchers. The ability to code is fast becoming an indispensable asset, as integral to bioscience as traditional laboratory skills.

## How Code Quality Improves Scientific Reproducibility

A main reason to write software is to define re-usable instructions for task automation (Narzt et al., 1998). However, the complexity of the code makes it prone to errors and can prevent usage by persons other than the author himself. This is a problem for the general scientific community, as the software is often the only way to reproduce the results of a study (Sandve et al., 2013). Hence, modern journals aim to enforce standards to software development, including software written and used by biological researchers (Smith et al., 2018). Here, we provide a brief overview of the standards utilized by `plotastic` that to ensure its reliability and reproducibility by the scientific community (Peng, 2011).

Modern software development is a long-term commitment of maintaining and improving

code after initial release (Boswell & Foucher, 2011). Hence, it is good practice to write the software such that it is scalable, maintainable and usable. Scalability or, to be precise, structural scalability means that the software can easily be expanded with new features without major modifications to its architecture (Bondi, 2000). This is achieved by writing the software in a modular fashion, where each module is responsible for a single function. Maintainability means that the software can easily be fixed from bugs and adapted to new requirements (Kazman et al., 2020). This is achieved by writing the code in a clear and readable manner, and by writing tests that ensure that the code works as expected (Boswell & Foucher, 2011). Usability is hard to define (Brooke, 1996), yet one can consider a software as usable if the commands have intuitive names and if the software’s manual, termed “documentation”, is up-to-date and easy to understand for new users with minimal coding experience. A software package that has not received an update for a long time (approx. one year) could be considered abandoned. Abandoned software is unlikely to be fully functional, since it relies on other software (dependencies) that has changed in functionality or introduce bugs that were not expected by the developers of all dependencies. Together, software that’s scalable, maintainable and usable requires continuous changes to its codebase. There are best practices that standardize the continuous change of the codebase, including version control, continuous integration (often referred to as CI), and software testing.

Version control is a system that records changes to the codebase line by line, allowing the documentation of the history of the codebase, including who made which changes and when. This is required to isolate new and experimental features into newer versions and away from the stable version that’s known to work. The most popular version control system is Git, which is considered the industry standard for software development (Chacon & Straub, 2024). Git can use GitHub.com as a platform to store and host codebases in the form of software repositories. GitHub’s most famous feature is called “pull request”. A pull request is a request from anyone registered on GitHub to include their changes to the codebase (as in “please pull this into your main code”). One could see pull requests as the identifying feature of the open source community, since it exposes the codebase to potentially thousands of independent developers, reaching a workforce that is impossible to achieve with closed source models used by paid software companies.

Continuous integration (CI) is a software development practice in which developers integrate code changes into a shared repository several times a day (Duvall et al., 2007). Each integration triggers the test suite, aiming to detect errors as soon as possible. The test suite includes building the software, setting up an environment for the software to run and then executing the programmed tests, ensuring that the software runs as a whole. Continuous integration is often used together with software branches. Branches are independent copies of the codebase that are meant to be merged back into the original code once the changes are finished. Since branches

accumulate multiple changes over time, this can lead to minor incompatibilities between the branches of all developers (integration conflicts), which is something that CI helps to prevent.

Continuous integration especially relies on a thorough software testing suite. Software testing is the practice of writing code that checks if the codebase works as expected (Myers et al., 2011). The main type of software testing is unit testing, which tests the smallest units of the codebase (functions and classes) in isolation (Listing 1).

**Listing 1:** Example of an arbitrary python function and its respective unit test function. The first function simply returns the number 5. The second function tests if the first function indeed returns the number 5. The test function is named with the prefix “test\_” and is placed in a file that ends with the suffix “\_test.py”. The test function is executed by the testing framework pytest. Note that code after “#” is considered a comment and won’t be executed.

```
1 # Define a function called "give_me_five" that returns the number 5
2 def give_me_five():
3     return 5
4 # Define a test function asserting that "give_me_five" returns 5
5 def test_give_me_five():
6     assert give_me_five() == 5
```

The quality of the software testing suite is measured by the code coverage, the precision of the tests, and the number of test-cases that are checked. The code coverage is the percentage of the codebase that is called by the testing functions, which should be as close to 100% as possible, although it does not measure how well the code is tested. The precision of the test is not a measurable quantity, but it represents if the tests truly checks if the code works as expected. The number of test-cases is the number of different scenarios that are checked by the testing functions, for example testing every possible option or combinations of options for functions that have multiple options. The most popular software testing framework for python is pytest, which is utilized by plotastic (Krekel et al., 2004).

Together, the standards of software development, including version control, continuous integration, and software testing, ensure that the software is scalable, maintainable, and usable. This is especially important for software that is used by the scientific community, as it ensures that the software is working as expected at defined versions years after publishing scientific results.

## Python as a Programming Language

Here, we provide a general overview of the python programming language, explaining terms like “*type*”, “*method*”, etc., in order to prepare readers without prior programming experience for the following chapters. We also describe the design principles of python to lay out the key concepts that differentiate python compared to other programming languages. A more detailed tutorial on python that’s specialized for bioscientists is found in Ekmekci et al. 2016

**Listing 2:** Example of readable python code. This one-line code returns the words (string) 'Hello, World!' when executed. The command is straightforward and easy to understand.

```
1 print("Hello, World!")
2 // Output: Hello, World!
```

Languages such as python are considered “*high-level*”, which means that it is designed to be easy to read and write, but also independent of hardware by hiding (“*abstracting*”) underlying details (*The Python Language Reference*, n.d.). A key principle of python is the emphasis on implementing a syntax that is concise and close to human language (Listing 2, Listing 3).

**Listing 3:** Example of less readable code written in the low-level programming language C. This code is doing exactly the same as the python code in Listing 2. The command is harder to understand because more steps are needed to access the same functionality, including the definition of a function

```
1 #include <stdio.h>
2 int main() {
3     printf("Hello, World!");
4     return 0;
5 }
6 // Output: Hello, World!
```

Furthermore, python is an *interpreted* language, which means that the code is executed line by line. This makes coding easier because the programmer can see the results of the code immediately after writing it, and error messages point to the exact line where the error occurred. This is in contrast to *compiled* languages, where the code has to be compiled into machine code before it can be executed. The advantage of compiled languages is that the code runs faster, because the machine code is optimized for the hardware.

Python automates tasks that would otherwise require an advanced understanding of computer hardware, like the need for manual allocation of memory space. This is achieved by using a technique called “*garbage collection*”, which automatically frees memory space that is no longer needed by the program. This is a feature that is not present in low-level programming languages like C or C++, that were designed to maximize control over hardware.

Another hallmark of python is its *dynamic typing system*. In python the type is inferred automatically during code execution (Listing 4). This is in contrast to *statically typed* languages like C, where the type of a variable has to be declared explicitly and cannot be changed during code execution (Listing 5) (*The Python Language Reference*, n.d.).

Dynamic typing makes python a very beginner-friendly language, since one does not have to keep track of the type of each variable. However, this also makes python a slower language, because the interpreter has to check the type of each variable during code execution. Also, developing code with dynamic typing systems is prone to introducing bugs (“*type errors*”), because it allows unexperienced developers to convert variables from one type to another without noticing, leading to unexpected behavior. Hence, larger python projects require disciplined

**Listing 4:** Example of dynamic typing in python. The variable “a” is assigned the value 5, which is of type integer. The variable “a” is then assigned the value “Hello, World!”, which is of type string. Python allows dynamic re-assignment of variables with different types. Note that code after “#” is considered a comment and won’t be executed.

```
1 a = 5 # Type integer
2 a = 5.0 # Type float
3 a = 'Hello, World!' # Type string
4 a = True # Type boolean
5 a = False # Type boolean
6 a = [1, 2, 3] # Type list of integers
7 a = {'name': 'Regina'} # Type dictionary
```

**Listing 5:** Example of static typing in C. The variable “a” is declared as an integer (int), and can only store integers. The variable “a” is then assigned the value 5, which is an integer. The variable “a” is then assigned the value ‘Hello, World!’, which is a string. This results in a compilation error, because the variable “a” can only store integers. Note that code after “//” is considered a comment and won’t be executed.

```
1 int a; // Declare type as integer
2 a = 5;
3 a = 'Hello, World!'; // Compilation error!
```

adherence to programming conventions. One such convention is *type hinting*, which is a way to explicitly note the type of a variable. Type hinting does not have an effect on the code, but it makes the code more readable and understandable for other developers, and allows for development environments to detect type errors before execution (Listing 6) (van Rossum et al., 2014).

**Listing 6:** Example of type hints used in python. Explicitly stating the type of the variable is optional and does not change the behavior of the code as shown in Listing 4.

```
1 a: int = 5
2 a: str = 'Hello, World!'
```

Python supports both functional and object-oriented programming paradigms. In functional programming, the code is written in a way that the program is a sequence of function calls, where each function call returns a value that is used in the next function call (Listing 7). This approach is useful when multiple actions have to be performed on the same data and the structure of the data is relatively simple, for example a string of a gene sequence.

When the data itself gains in complexity, for example when storing not just the gene sequence, but also the promotor sequence, an object-oriented approach is more suitable (Listing 8). Object-oriented programming is a programming paradigm that uses objects and classes. An object is a collection of both data and functions, and a class is a blueprint for creating objects. The data of an object is stored as attributes. Functions that are associated with an object are called methods.

**Listing 7:** Example of functional programming in Python. The code defines a function called “find\_restriction\_site” that finds the position of a restriction site in a gene. The function “cut” uses the function “find\_restriction\_site” to cut the gene at the restriction site.

```

1 def find_restriction_site(gene: str):
2     return gene.find('GCGC')
3
4 def cut(gene: str):
5     position = find_restriction_site(gene)
6     return gene[:position]
7
8 gene1 = 'TGAGCTGAGCTGATGCGCTATATTAGGCG'
9 gene1_cut = cut(gene1)
10 print(gene1_cut)
11 # Output: GCGCTATATTAGGCG

```

**Listing 8:** Example of object oriented programming in python. The class is called “Gene” and has four methods, “\_\_init\_\_”, “find\_promotor”, “find\_restriction\_site” and “cut”. The method “\_\_init\_\_” is called when creating (“initializing”) an object, which fills the object with user-defined data. The parameter “self” is used to reference the object itself internally. “find\_promotor” is a method that finds the position of the promotor in the gene and is called during object initialization.

```

1 class Gene:
2     def __init__(self, sequence: str):
3         self.sequence: str = sequence # Save sequence as attribute
4         self.promotor: str = self.find_promotor()
5     def find_promotor(self):
6         return self.sequence.find('TATA')
7     def find_restriction_site(self):
8         return self.sequence.find('GCGC')
9     def cut(self):
10        position = self.find_restriction_site()
11        return self.sequence[:position]
12
13 gene1 = Gene(sequence='TGAGCTGAGCTGATGCGCTATATTAGGCG') # Create object
14 gene1_cut = gene1.cut() # Call the method cut
15 print(gene1_cut) # Show result
16 # Output: GCGCTATATTAGGCG

```

A major benefit of using an object oriented versus a functional approach is that the data itself is programmable, enabling the programmer to define the behavior of the data itself through methods. This is achieved by using the keyword “self” to reference the object itself inside the class. For example, one could extend the class “Gene” with a method that finds the promotor of the gene and stores it as an attribute (Listing 8).

When designing software, both functional and object oriented programming can be used together, where object oriented programming is often used to design the program’s overall architecture, and functional programming is used to implement the algorithms of the program’s features. This allows for scalability of the software, as every single class is extended through the

addition of new methods. Furthermore, classes can be expanded in their functionalities through inheritance (Listing 9).

**Listing 9:** Example of inheritance in python. The class “mRNA” inherits from the class “Gene”. The class “mRNA” has two methods, “\_\_init\_\_” and “find\_stopcodon”. The method “find\_stopcodon” finds the position of stop codons.

```
1 # Define a class called mRNA inheriting from the class Gene
2 class mRNA(Gene):
3     def __init__(self, sequence: str):
4         super().__init__(sequence) # Get attributes from parent class
5         self.sequence.replace('T', 'U') # Replace thymine with uracil
6     def find_stopcodons(self):
7         return self.sequence.find('UGA')
8
9 mRNA1 = mRNA(sequence='TGAGCTGAGCTGATGCGCTATTTAGGGC') # Create object
10 print(mRNA1.find_stopcodons()) # Call the method translate
11 # Output: [0, 5, 10]
```

Inheritance is a feature of object-oriented programming that allows a class to access every attribute and method of a parent class. For example, one could extend the class “Gene” with a class “mRNA”, by writing a class “mRNA” that inherits from the class “Gene”.

Together, python is not just beginner-friendly, but also well respected for its ease in development, which is why it is widely used in professional settings for web development, data analysis, machine learning, biosciences and more (Ekmekci et al., 2016).

## Data Science with Python

the ease of use has made python a very popular language (Rayhan & Gross, 2023)

Like any other programming language, python alone does not provide specialized tools like those used for data analysis (*The Python Language Reference*, n.d.). However, python was designed to be extended by packages developed by its users. A python package consists of multiple python modules, where each module is a text-file with a .py ending containing python code. Famous examples of such packages are pytorch and tensorflow, that are used to build models of artificial intelligence, including ChatGPT (Paszke et al., 2019; Abadi et al., 2016; Radford et al., 2019). Here, we outlay the most important packages used for plotastic.

Interactive Python - Jupyter

Python overcame the issues of interpreted language by utilizing Code written in C numpy:

- Acceleration, - SIMD instructions

Tabular operations - pandas

Data visualization - matplotlib - seaborn

Inferential Statistics - pingouin

AI: - pytorch and tensorflow - example: VGG19 is just a few lines of code (??) asdfdf

## Aims

This project defines these aims:

- Characterise the interaction between myeloma cells and mesenchymal stromal cells
- Aim 2
- Aim 3

# Chapter 1: Modelling Myeloma Dissemination *in vitro*

## Abstract

lorem ipsum

## Introduction

lorem ipsum dolor sit amet

## Methods

lorem ipsum

## Results

lorem ipsum

## Discussion

lorem ipsum

## Supplementary Figures and Methods

### Keep it Together: Describing Myeloma Dissemination *in vitro* with hMSC-Interacting Subpopulations and their Aggregation/Detachment Dynamics

Martin Kuric<sup>1</sup>, Susanne Beck<sup>2</sup>, Doris Schneider<sup>1</sup>, Wyonna Rindt<sup>3</sup>, Marietheres Evers<sup>4</sup>, Jutta Meißner-Weigl<sup>1</sup>, Sabine Zeck<sup>1</sup>, Melanie Krug<sup>1</sup>, Marietta Herrmann<sup>5</sup>, Tanja Nicole Hartmann<sup>6</sup>, Ellen Leich<sup>4</sup>, Maximilian Rudert<sup>7</sup>, Denitsa Docheva<sup>1</sup>, Anja Seckinger<sup>8</sup>, Dirk Hose<sup>8</sup>, Franziska Jundt<sup>3</sup>, Regina Ebert<sup>1</sup>

<sup>1</sup>University of Würzburg, Department of Musculoskeletal Tissue Regeneration, Würzburg, Germany

<sup>2</sup>University Hospital Heidelberg, Institute of Pathology, Heidelberg, Germany

<sup>3</sup>University Hospital Würzburg, Department of Internal Medicine II, Würzburg, Germany

<sup>4</sup>University of Würzburg, Institute of Pathology, Comprehensive Cancer Center Mainfranken, Würzburg, Germany

<sup>5</sup>University Hospital Würzburg, IZKF Research Group Tissue Regeneration in Musculoskeletal Diseases, Würzburg, Germany

<sup>6</sup>University of Freiburg, Department of Internal Medicine I, Faculty of Medicine and Medical Center, Freiburg, Germany

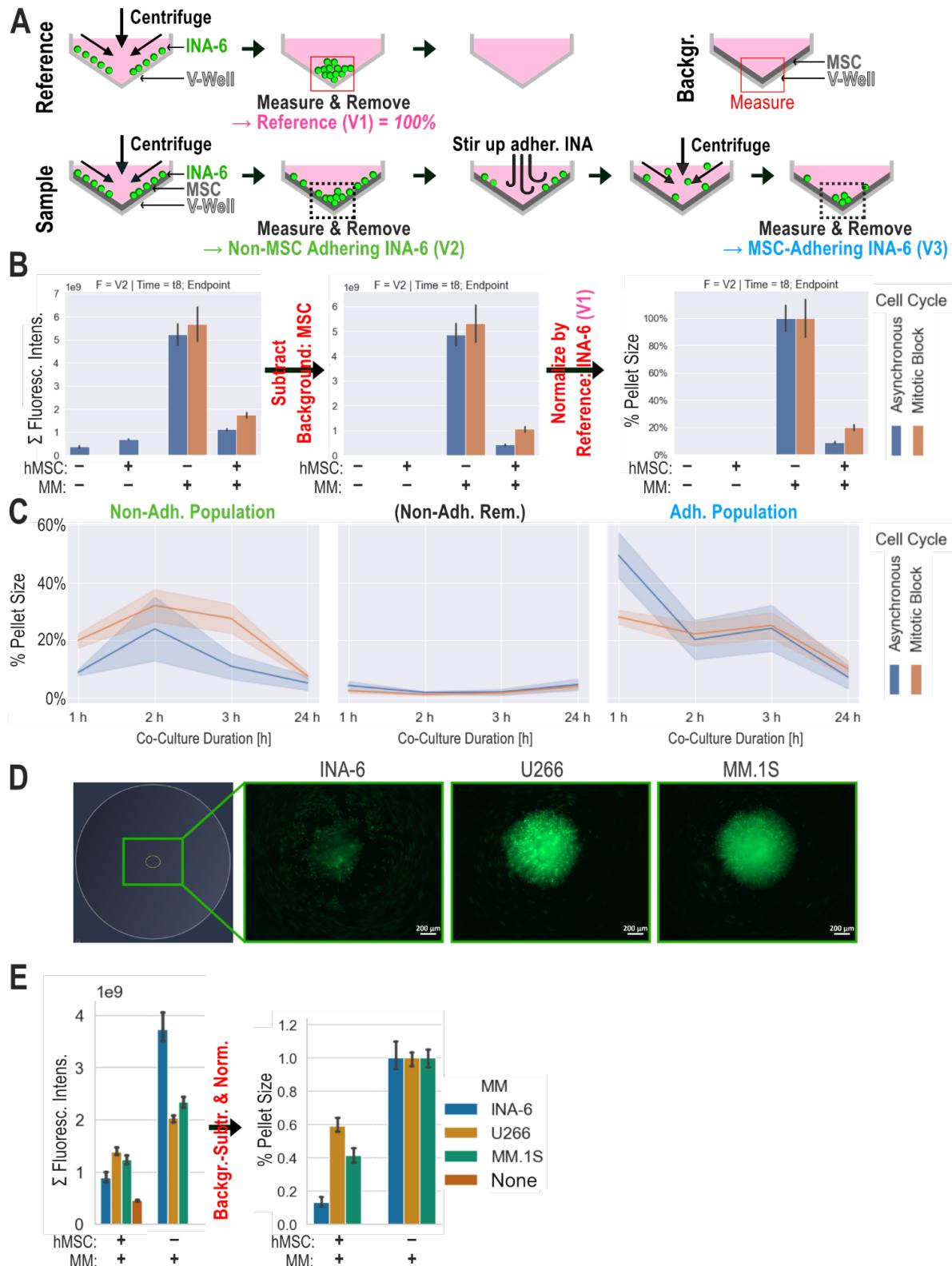
<sup>7</sup>University of Würzburg, Orthopedic Department, Clinic König-Ludwig-Haus, Würzburg, Germany

<sup>8</sup>Vrije Universiteit Brussel, Department of Hematology and Immunology, Jette, Belgium

**Tab. S1:** List of hMSC donors, myeloma cell lines, and their mycoplasma test status. If no unique donors were available, hMSC donors were used twice for the same experiment at different passages. WPSC: Well plate sandwich centrifugation.

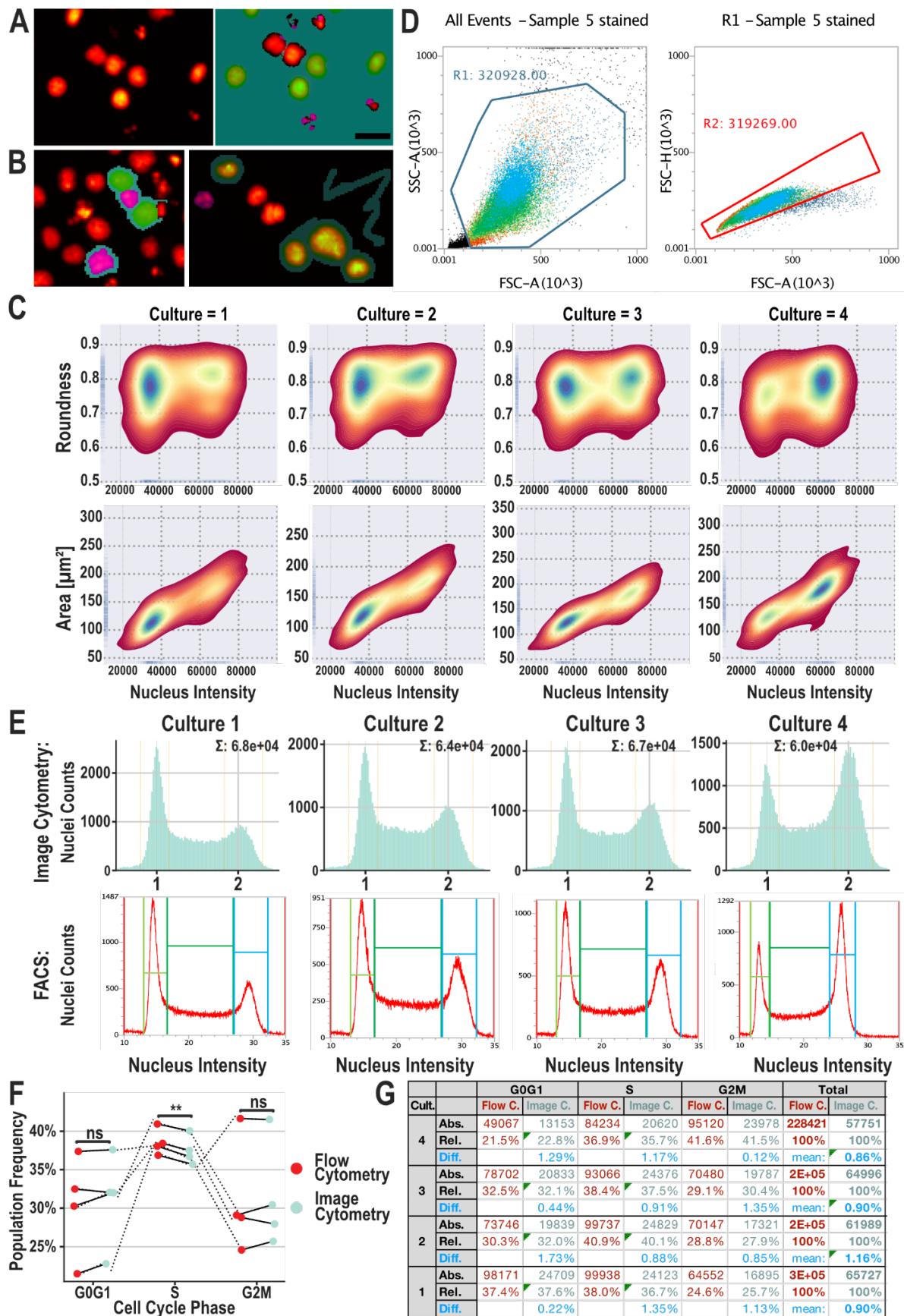
Cell Type	Donor / Line	Donor Ages	Donor Sex	Date of negative Mycoplasma test	Experiment(s)	Figures
Myeloma Cell Line	<b>INA-6</b>	80	m	09.02.22	All	All
	<b>U266</b>			10.10.22	- Validation of V-Well Adhesion Assay	S1E
	<b>MM1.S</b>			24.02.22		
hMSC	<b>1639</b>	49	m	not tested	- Validation of V-Well Adhesion Assay - Time-lapse: INA-6 on dispersed hMSC	S1E 1D; 2[A-E]
	<b>1571</b>	72	m	not tested	- Saturation of hMSCs	1[A-B]
	<b>1573</b>	47	m	not tested		
	<b>1578</b>	82	m	not tested		
	<b>1842</b>	63	m	not tested	- INA-6 Viability dep. on time and hMSC adhesion surface (INA not washed off)	1E right
	<b>1843</b>	60	m	not tested		
	<b>1537</b>	77	f	not tested		
	<b>1794</b>	82	m	not tested		
	<b>1779</b>	61	m	not tested	- INA-6 Viability dep. on time and hMSC adhesion surface (INA washed off)	1[C, E left]
	<b>1849</b>	69	m	not tested		
	<b>1854</b>	80	f	not tested		
	<b>1605</b>	71	f	not tested	- Time-lapse: INA-6 on dispersed hMSC	1D; 2[A-E]
	<b>1650</b>	57	m	not tested		
	<b>1859</b>	64	f	not tested	- Time-lapse: INA-6 on confluent hMSC	2[G-I]
	<b>1863</b>	79	f	not tested		
	<b>1861</b>	52	f	not tested		
	<b>1818</b>	81	f	not tested	- Cell Cycle Profiling after V-well assay	3C
	<b>1824</b>	82	f	not tested	(Donor measured twice, different passages) - V-well adhesion assay of mitotically blocked INA-6 followed by Cell Cycle Profiling after V-well assay	3[B,C]
	<b>1827</b>	56	m	not tested	- V-well adhesion assay of mitotically blocked INA-6 followed by Cell Cycle Profiling after V-well assay	

				assay	
<b>1501</b>	59	m	not tested	- INA-6 AI-assisted count during WPSC (INA-6 stained with celltracker green)	4B
<b>1643</b>	75	f	not tested		
<b>1718</b>	67	m	not tested		
<b>1720</b>	58	m	not tested		
<b>1653</b>	65	m	not tested		
<b>1591</b>	78	m	not tested	- WPSC (MACS) followed by RNAseq, Metascape analysis and qPCR validation - WPSC (Wash) followed by qPCR-Validation and Luminescent Viability assays	4[A,C,D,E] ; 5[A-C] 4[C-E], 4F
<b>1654</b>	74	m	not tested	- WPSC (MACS) followed by RNAseq, Metascape analysis and qPCR validation - WPSC (Wash) followed by qPCR-Validation and Luminescent Viability assays	4[A,C,D,E] ; 5[A-C] 4[C-E], 4F
<b>1655</b>	78	f	not tested	- WPSC (MACS) followed by RNAseq, Metascape analysis and qPCR validation	4[A,C,D,E] ; 5[A-C]
<b>1668</b>	80	f	not tested		
<b>1670</b>	66	f	not tested		
<b>1701</b>	81	m	not tested	- WPSC (Wash) followed by qPCR-Validation and Luminescent Viability assays	4[C-E], 4F
<b>1702</b>	79	f	not tested		
<b>1600</b>	77	m	not tested		
<b>1681</b>	56	m	not tested	- WPSC (Wash) followed by Luminescent Viability assays	4F
<b>1672</b>	65	m	not tested	- WPSC (Wash) followed by qPCR-Validation	4[C-E]

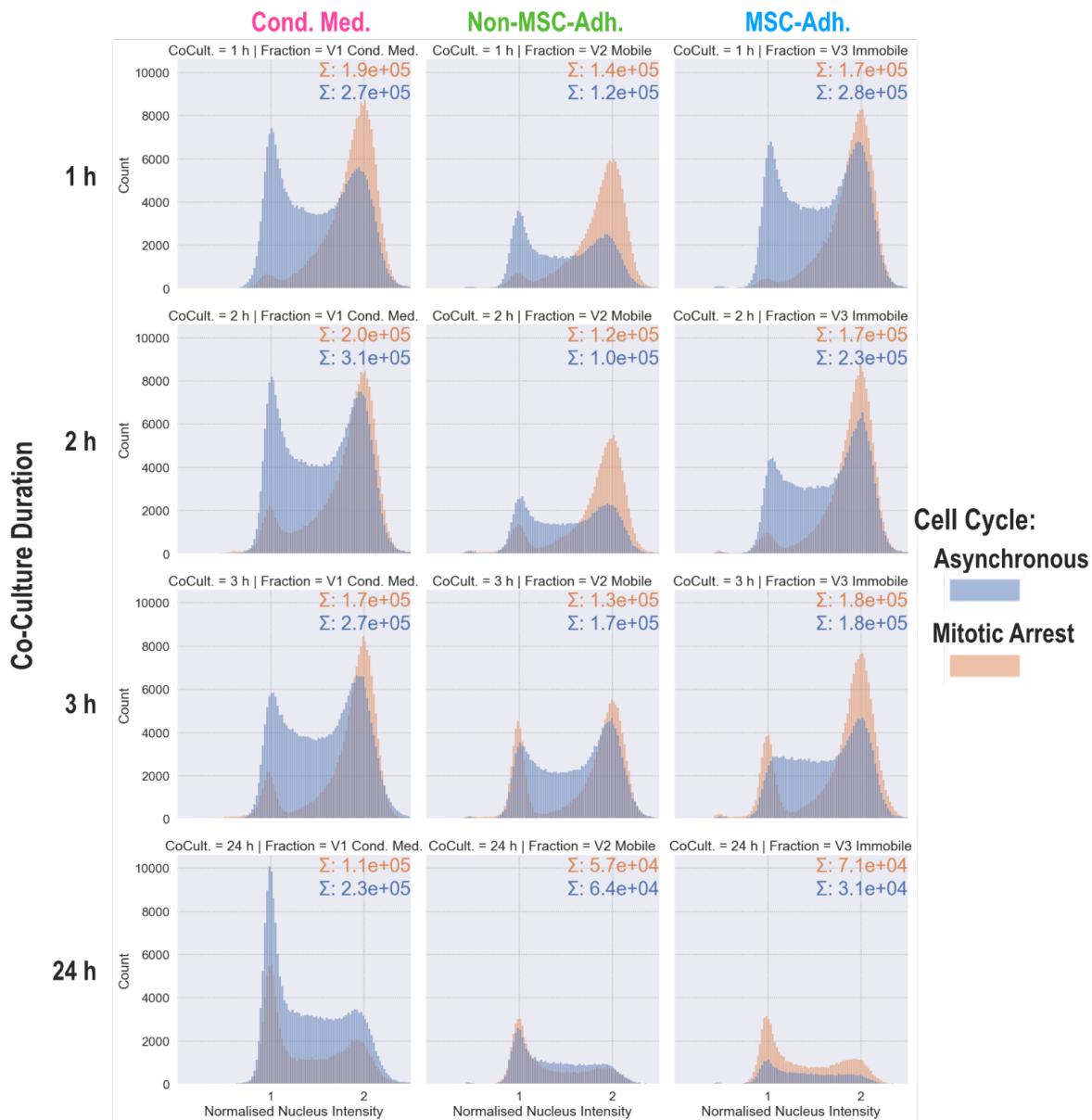
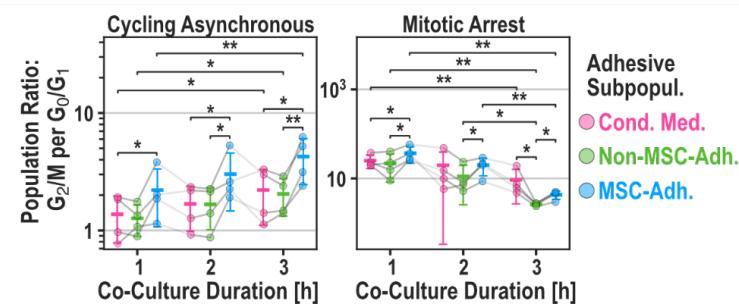


**Fig. S1:** Principle and quantification of the V-well adhesion assay of fluorescently labeled myeloma cells adapted by Weetall et al. 2001. **A:** Sample: Subsequent rounds of centrifugation and removal of cell pellet yielded the size of adhesive subpopulations. Fluorescently stained INA-6 cells were added to an hMSC monolayer. Non-adherent INA-6 cells (V2) were pelleted in the well-tip. Pellets were quantified by fluorescence brightness and isolated by pipetting. Immobile INA-6 cells (V3) were manually detached by forceful pipetting.

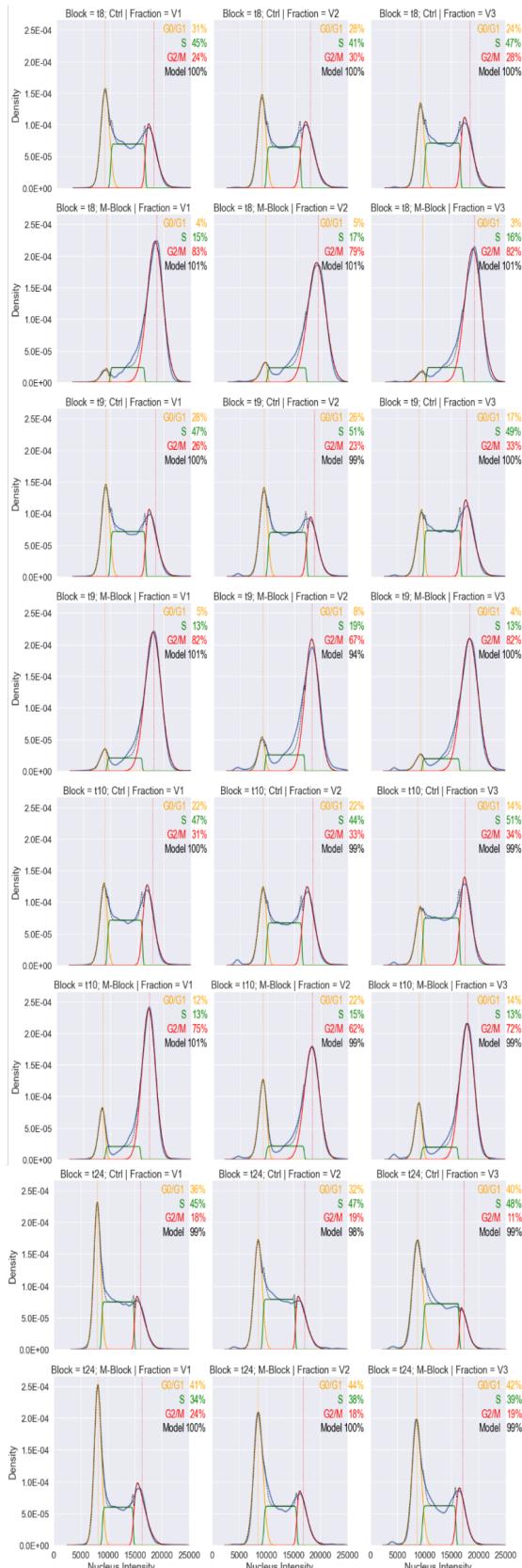
Reference: Omitting adhesive hMSC-layer yielded ~100% non-adherent cells (V1) after the first centrifugation step; Background: hMSC monolayer was used as background signal. **B:** Calculation of the population size relative to total cells starting with pellet intensity. The shown example is the pellet gained by centrifuging mobile subpopulation (V2) after 1 h of co-culture. (see Fig. 2 for context): Intensity values from pellet images were summarized. After subtracting the unlabeled hMSC signal and normalization by a full-size pellet (reference), the resulting values represented the fraction of the adhesive subpopulation. **C:** One of three biological replicates summarized in Fig. 2. Line range shows the standard deviation of four technical replicates. Non.Adh. Rem.: Fluorescence signal after removal of V2. **D:** Example images of myeloma cell lines (INA-6, U266, MM.1S) pelleted in the tip of V-wells. The leftmost image shows the recorded area in a complete V-well. Scale bar = 200  $\mu$ m. **E:** Results from (D) comparing adhesion strength of three myeloma cell lines to hMSC. Error bars represent technical deviation. MM=Multiple Myeloma.



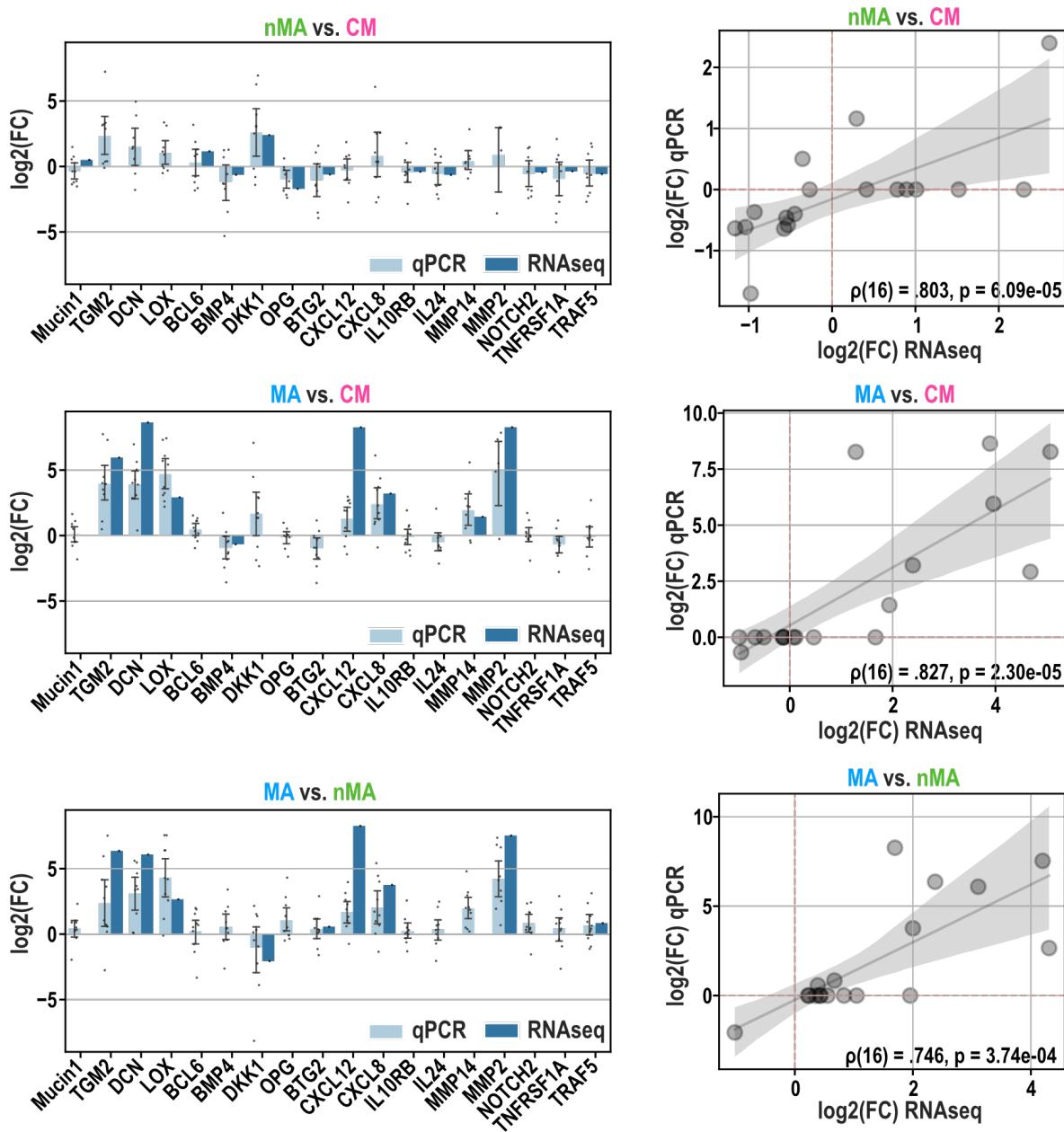
**Fig. S2:** Validation of image cytometric analysis of cell cycle in four INA-6 cultures **A:** Left: Example image cytometric scan: INA-6 cells were stained with Hoechst33342 and scanned by automated fluorescence microscopy. Right: The image was segmented using a convolutional neural network (ZEISS ZEN intellesis) trained to discern healthy nuclei (green) from fragmented ones (magenta). Doublets are excluded by setting an area- and roundness threshold. Scale bar: 20  $\mu$ m. **B:** Two example images from the training set. **C:** Quality of image cytometric data was ensured by plotting the distribution of nuclei brightnesses vs. the distribution of both nuclei-roundnesses and nuclei-areas. Nuclei with double fluorescence intensity have the same roundness while their area increases, as expected from a cell in G2 phase. **D:** The same samples from (C) were also measured with flow cytometry. Representative example of gating strategy: Left: Dead cells were excluded by setting a minimum threshold for side-scattering (SSC-A). Right: Doublets were excluded by setting a maximum threshold for forward scatter area (FSC-A) (sample “5” represents culture “4” in this figure). **E:** Cell cycle profiles of four independent INA-6 cultures were measured by both image cytometry (top) and flow cytometry (bottom). For both methods, frequencies of G0/G1, S, and G2M were summed up by setting fluorescence intensity thresholds. **F:** Image cytometry yields the same frequencies for G0/G1, S, and G2M when compared to flow cytometry. RM-ANOVA showed that the method has no significant effect on the frequencies of cell cycle populations [ $F(1,3)=1.421$ ,  $p\text{-unc}=.32$ ]. **G:** Results from (F) in tabular form. On average, frequencies for G0/G1, S, and G2M measured by Image cytometry differ by 0.95 percent points compared to flow cytometry measurement. Cult.: Culture; C.: Image cytometry; Abs.: Absolute cell count; Rel.: Relative cell count; Diff.: Difference between relative cell counts determined by flow cytometry and image cytometry.

**A****B**

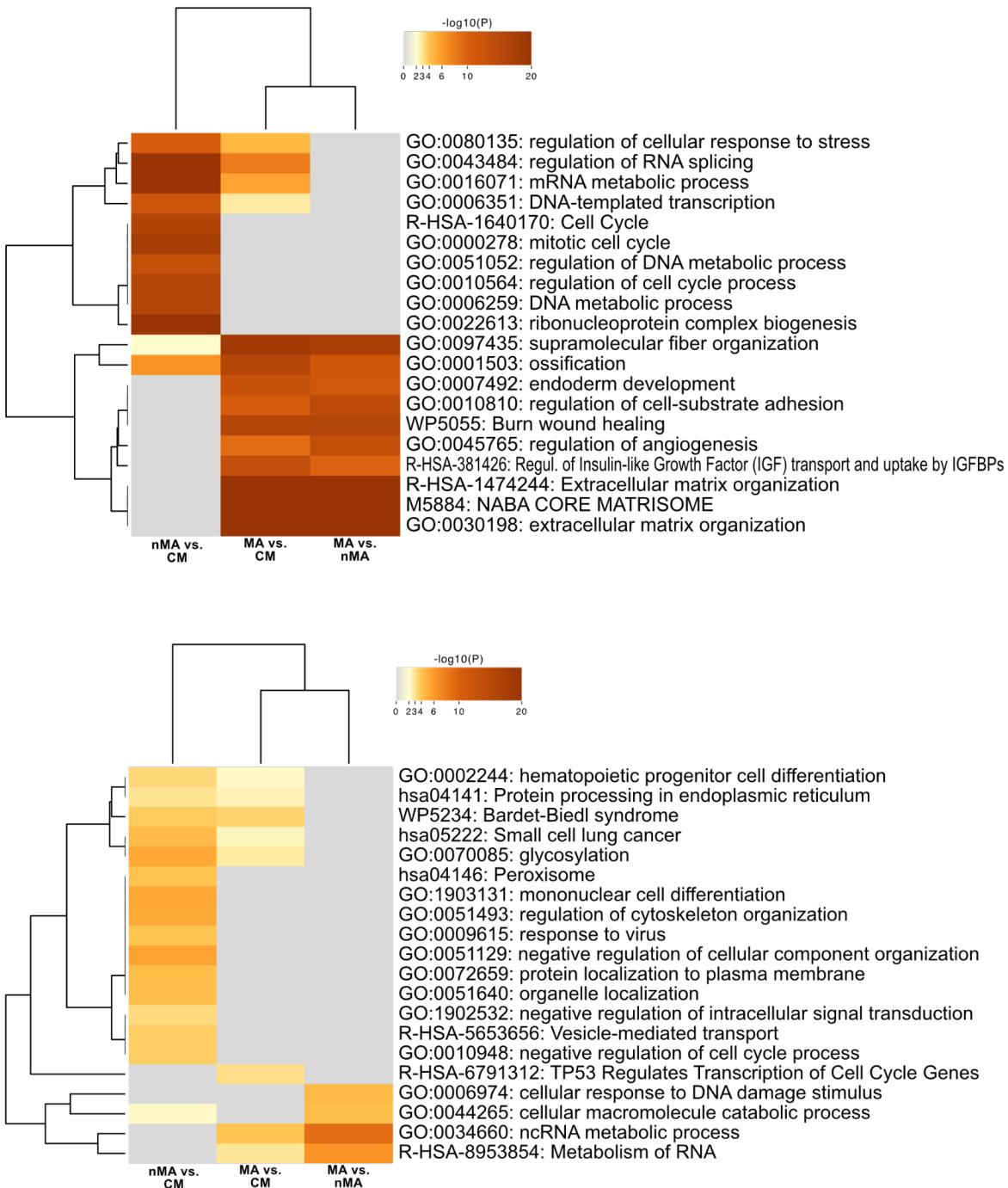
**Fig. S3:** Cell cycle analysis of INA-6 pellets gained from V-Well Adhesion assay (Fig. 3). **A:** Cell cycle profiles of MSC-adhering subpopulations. INA-6 cells were synchronized by double thymidine block followed by nocodazole. Cell cycle was released directly before addition to hMSCs. Histograms were normalized and summed up across all biological replicates ( $n=4$ ). Technical replicates (3) were pooled prior to cell cycle profiling. CoCult. = Co-culture duration. Fraction = Adhesion subpopulations. **B:** Similar figure to Fig. 3C displaying ratio of INA-6 populations ( $G_2/M$  to  $G_0/G_1$ ). **Statistics:** Paired t-test (B).



**Fig. S4:** Representative (one of the four independent sample sets as seen in Fig. S3) curve fitting analysis of cell cycle profiles generated by Image Cytometry. t8, t9, t10, and t24 refer to 1, 2, 3, and 24 hours after the addition of INA-6 cells to hMSCs.



**Fig. S5:** Correlation of RNAseq with qPCR **Left:** Validation of RNAseq results (Fig. 3) with qPCR showing the  $\log_2(\text{foldchange expression})$  of 18 genes. For qPCR, Datapoints each represent one biological replicate ( $n=10$ ), which is the mean of technical replicates ( $n=3$ ). Bar height represents mean of biological replicates, error bars show standard deviation of biological replicates. **Right:** Correlation between qPCR and RNAseq in terms of  $\log_2(\text{mean foldchange expression per gene})$ . Each dot represents one gene shown in the barplot to the left. Genes measured with qPCR that showed no differential expression in RNAseq were set to have a  $\log_2(\text{FC}) = 0$ . Shaded area shows the confidence interval of linear regression. Correlation coefficient ( $\rho$ ) was calculated using Spearman's rank. N = 18 genes. FC = fold change expression.



**Fig. S6:** Functional enrichment analysis by Metascape using genes that are differentially expressed between MSC-interacting subpopulations. **Top:** Upregulated genes. **Bottom:** Downregulated genes.

**Tab. S2:** Adhesion genes (from Fig. 6A) filtered by association with patient survival ( $p < 0.05$ ) and categorized by a continuous downregulation across disease progression. The full table including non-significant associations is found in the supplementary data. Bone Marrow Plasma Cell (BMPC), Monoclonal Gammopathy of Undetermined Significance (MGUS), Smoldering Multiple Myeloma (sMM), Multiple Myeloma (MM), Multiple Myeloma Relapse (MMR). p-adj. = adjusted p-values (Benj.-Hoch.).

Regulation during disease progression	Gene	Ensemble ID	Progression on Free / Overall Survival	Better Prognosis with high/low expression	Association of expression with survival	
					[p-unc]	[p-adj]
Not downregulated or overall low expression	CCDC80	ENSG00000091986	Prog. Free	high	2.04E-03	1.58E-02
	CCN2	ENSG00000118523	Overall	high	2.89E-03	2.43E-02
	CCNE2	ENSG00000175305	Prog. Free	low	1.21E-02	4.62E-02
			Overall	low	5.34E-04	8.64E-03
	COL4A1	ENSG00000187498	Overall	high	9.47E-03	3.99E-02
	COL4A2	ENSG00000134871	Prog. Free	high	1.24E-02	4.62E-02
	F3	ENSG00000117525	Overall	high	9.18E-03	3.99E-02
	HTRA1	ENSG00000166033	Prog. Free	high	1.20E-02	4.62E-02
	IGFBP7	ENSG00000163453	Prog. Free	low	9.53E-03	4.38E-02
	MMP2	ENSG00000087245	Prog. Free	high	2.29E-05	2.32E-03
	OSMR	ENSG00000145623	Prog. Free	high	5.67E-04	7.15E-03
			Overall	high	1.29E-02	4.64E-02
Continuously downregulated (PC > MGUS > sMM > MM > MMR)	SERPINH1	ENSG00000149257	Prog. Free	low	1.83E-03	1.58E-02
			Overall	low	4.40E-03	2.61E-02
	ACTN1	ENSG00000072110	Overall	high	7.73E-03	3.55E-02
	AEBP1	ENSG00000106624	Prog. Free	high	1.08E-02	4.62E-02
	AXL	ENSG00000167601	Prog. Free	high	1.50E-03	1.51E-02
			Overall	high	3.64E-05	1.84E-03
	COL1A1	ENSG00000108821	Prog. Free	high	3.03E-04	4.37E-03
			Overall	high	5.93E-04	8.64E-03
	COL3A1	ENSG00000168542	Overall	high	1.08E-02	4.29E-02
	COL6A1	ENSG00000142156	Prog. Free	high	1.20E-02	4.62E-02
			Overall	high	1.10E-02	4.29E-02
	CXCL12	ENSG00000107562	Prog. Free	high	1.16E-04	2.93E-03
			Overall	high	6.48E-04	8.64E-03
	CYP1B1	ENSG00000138061	Prog. Free	high	8.64E-03	4.17E-02
			Overall	high	6.84E-04	8.64E-03
	DCN	ENSG00000011465	Prog. Free	high	4.83E-03	3.05E-02
			Overall	high	2.47E-04	8.33E-03
	FBLN1	ENSG00000077942	Prog. Free	high	2.68E-03	1.93E-02
			Overall	high	3.73E-03	2.61E-02
	GNB3	ENSG00000111664	Prog. Free	high	3.75E-03	2.52E-02
			Overall	high	5.73E-03	3.05E-02
	IGFBP4	ENSG00000141753	Prog. Free	high	8.68E-03	4.17E-02
			Overall	high	7.09E-03	3.41E-02
	ITGAX	ENSG00000140678	Prog. Free	high	6.72E-03	3.60E-02
			Overall	high	3.12E-03	2.43E-02
	LAMB2	ENSG00000172037	Overall	high	1.35E-03	1.39E-02
	LRP1	ENSG00000123384	Prog. Free	high	6.46E-03	3.60E-02
			Overall	high	4.34E-04	8.64E-03
	LTBP2	ENSG00000119681	Prog. Free	high	9.03E-05	2.93E-03
			Overall	high	1.17E-02	4.36E-02
	MAP3K8	ENSG00000107968	Prog. Free	high	9.58E-04	1.08E-02
	MFAP5	ENSG00000197614	Prog. Free	high	2.43E-04	4.09E-03
			Overall	high	4.27E-03	2.61E-02
	MMP14	ENSG00000157227	Prog. Free	high	6.93E-05	2.93E-03
			Overall	high	6.69E-03	3.38E-02
	MYL9	ENSG00000101335	Prog. Free	high	1.46E-04	2.95E-03
			Overall	high	1.56E-05	1.57E-03
	NRP1	ENSG00000099250	Prog. Free	high	1.89E-03	1.58E-02
			Overall	high	2.21E-03	2.03E-02
	TGFBI	ENSG00000120708	Overall	high	4.30E-03	2.61E-02
	TNC	ENSG00000041982	Prog. Free	high	1.28E-02	4.62E-02
			Overall	high	4.75E-03	2.67E-02
	TPM1	ENSG00000140416	Overall	high	1.37E-03	1.39E-02
	TUBA1A	ENSG00000167552	Prog. Free	low	6.78E-03	3.60E-02

## Supplementary Materials and Methods

### **Isolation and Culturing of Primary Human Bone Marrow-Derived Mesenchymal Stromal Cells**

Primary human MSCs were obtained from the femoral head of patients (S. Tab. 1) undergoing elective hip arthroplasty. Material was collected with the informed consent of all patients and the procedure was approved by the local Ethics Committee of the University of Würzburg (186/18). In brief, bone marrow was washed with MSC-Medium [Dulbecco's modified Eagle's medium (DMEM/F12) (Thermo Fisher Scientific, Darmstadt, Germany) supplemented with 10% Fetal Calf Serum (FCS) (Bio&Sell GmbH, Feucht, Germany (Fernandez-Rebollo et al., 2017), 100 U/ml penicillin, 0.1 mg/ml streptomycin, 50 µg/ml ascorbate and 100 nmol/l sodium selenite (both Sigma-Aldrich GmbH, Munich, Germany)) and centrifuged at 250 g for 5 min. The pellet was washed four times with MSC-medium and resulting supernatants containing released cells were collected. Cells were pelleted and cultured at a density of  $1 \times 10^9$  cells per 175 cm<sup>2</sup> culture flask. After two days non-attached cells were washed away and adherent ones were cultivated in MSC-Medium until confluence. Then, they were either frozen in liquid nitrogen or directly utilized for experiments. hMSC cultures were sustained for a maximum of two passages. All cells were cultured at 37 °C and at 5% CO<sub>2</sub>.

### **Culturing of Myeloma Cell Lines**

The plasmacytoma cell line INA-6 [RRID:CVCL\_5209; DSMZ, Braunschweig, Germany, authenticated by DSMZ 2014 (see supplemental); (Burger et al., 2001; Gramatzki et al., 1994) was cultivated in RPMI 1640 medium (Life Technologies GmbH) supplemented with 20 % (v/v) FCS, 100 µg/ml gentamicin, 2 mmol/l L-glutamine (both Life Technologies GmbH), 1 mmol/l sodium pyruvate, 100 nmol/l sodium selenite (both Sigma Aldrich GmbH, Schnelldorf, Germany) and 2 ng/ml recombinant human interleukin-6 (IL-6; Miltenyi Biotec, Bergisch Gladbach, Germany). INA-6 were passaged three times per week by diluting them to  $1 \times 10^5$ ,  $2 \times 10^5$ , or  $4 \times 10^5$  cells/mL for 3, 2 and 1 days of culturing, respectively. MM.1S (RRID:CVCL\_8792) (Greenstein et al., 2003), and U266 cells (CVCL\_0566), (Nilsson et al., 1970) were propagated and cultivated in RPMI1640 medium comprising 10 % (v/v) FCS, 100 U/ml penicillin, 100 µg/ml streptomycin, 2 mmol/l L-glutamine, and 1 mmol/l sodium pyruvate. All cells were cultured at 37 °C and at 5% CO<sub>2</sub>.

**Co-Culturing of Primary hMSCs and INA-6 and MSC-Conditioning of Medium**

For each co-culture, hMSCs were seeded out 24 h prior to INA-6 addition to generate MSC-conditioned medium (CM). CM from different donors was collected separately and used immediately when adding INA-6. To ensure that CM was free of hMSCs, it was strained (40 µm) and centrifuged for 15 minutes at 250 g. INA-6 cells were washed with PBS (5 min, 1200 rpm), resuspended in MSC-medium and added to hMSCs such that co-culture comprised 33% (v/v) of CM gathered directly from the respective hMSC-donor. Co-cultures did not contain IL-6 (Chatterjee et al., 2002).

**Collagen I Coating**

Collagen I solution (isolated from rat tail, Corning, NY, USA) was diluted 1:2 (75 ng/mL) in acetic acid (0.02 N), applied to 96 well plates (30 µL in each well) and incubated for 2 h at room temperature. Acetic acid was removed and wells were washed once with 100 µL of PBS. Coated plates were stored dry at 4 °C.

**Fluorescent Staining of Cells**

For each live staining, cells were strained (70 µm) to remove clumps and washed (5 min, 250 g) once with the respective media (without FCS) and then resuspended in staining reagents.

For CellTracker™ Green CMFDA Dye and CellTracker™ Deep Red Dye (Thermo Fisher Scientific) staining, 1 mL staining solution for a maximum of  $1 \times 10^6$  cells was prepared. Staining was done at RT for 15 min using 5 µM CMFDA (5-Chlormethyl-fluoresceindiacetat) and 5 min of 1-2 µM DeepRed. To reduce background, stained cells were pelleted, resuspended in cell medium (containing FCS), incubated for 30 min (37 °C, 5% CO<sub>2</sub>), washed in cell medium, resuspended in 100 µL - 1 mL and counted.

For PKH26 staining (Sigma Aldrich), a maximum of  $1 \times 10^4$  cells was resuspended in 500 µL diluent C before swiftly adding 500 µL of staining solution (1 µL diluted in 500 µL diluent C) and incubating cells for 5 min at RT. The staining reaction was stopped by adding 1 mL of FCS-containing medium and adding 3 mL of FCS-free medium. Cells were washed with 10 mL of FCS-containing medium, resuspended in 100 µL - 1 mL cell medium, and counted.

For Calcein-AM (Calcein-O,O'-diacetat-tetrakis-(acetoxymethyl)-ester) (Thermo Fisher Scientific) staining, end concentrations of 0.5 µM were used. 12.5 µL of diluted stock solution (2.5 µM) was carefully added to 50 µL of the co-culture and incubated for 10 minutes at 37 °C.

For Hoechst33342 staining, cells were washed once with PBS, resuspended in a maximum of 500 µL of PBS, and fixed with 5 mL of ice-cold ethanol (70% v/v) by vigorously pipetting up and down to dissociate aggregates. Cells were washed once with PBS and stained with 2.5 µg/mL Hoechst33342 (Thermo Fisher Scientific) diluted in PBS for 1 h at 37 °C.

### **Automated Fluorescence Microscopy**

To remove clumps for microscopic applications, we cultured cells in 40 µm strained FCS. To reduce background fluorescence and phototoxicity, we used phenol-red free versions of the respective medium, if available.

All microscopy equipment was acquired from ZEISS. The microscope was an Axio Observer 7 with confocal Apotome.2 equipped with a motorized reflector revolver and motorized scanning table (130x100 mm). The microscope was mounted on an Antivibrations-Set [Axio Observer (D)] with two antivibration carrier plates, each equipped with two vibration dampening feet. The light source was a microLED 2 for transmission light and (for fluorescence) Colibri 7 (R[G/Y]B-UV) for five channels of incident light (385, 475, 555, 590, 630 nm). For excitation (EX) and emission (EM) light filtering and beam splitting (BS) we used the following reflectors: 96 HE BFP shift free (E) (EX: 390/40, BS: 420, EM: 450/40), 43 HE Cy 3 shift free (E) (EX: 550/25, BS: 570, EM: 605/70), 38 HE eGFP shift free (E) (EX: 470/40, BS: 495, EM: 525/50) and 90 HE LED (E) (EX: 385, 475, 555 und 630 nm, BS: 405 + 493 + 575 + 653, EM: 425/30 + 514/30 + 592/30 + 709/100). We used the black and white camera Axiocam 506 mono (D) and if not stated otherwise, 2x2 binning was used for fluorescence imaging. For mosaic acquisitions (“tiles”) we used a tiling overlap of 8-10% and image tiles were not stitched. Images were magnified 5x and 10x (Fluar 5x/0.25 M27 and EC Plan-Neofluar 10x/0.3 Ph1 M27).

### **Cell Viability and Apoptosis Assay**

To examine cell viability and apoptosis, cells were seeded in a 96-well plate ( $1 \times 10^4$  cells per well) to be measured inside culture well after respective incubation time immediately. ATP-amount and Caspase 3/7 activity were used as a proxy for viability and apoptosis rates, respectively. They were assessed using the CellTiter-Glo Luminescent Cell Viability Assay and the Caspase-Glo 3/7 Assay, respectively (Promega GmbH, Mannheim, Germany), according to the manufacturer's instructions.

Luminescence was measured with an Orion II Luminometer (Berthold Detection Systems, Pforzheim, Germany).

### **Microscopic Characterization of MSC Saturation**

For saturating hMSC with INA-6, hMSCs were stained with CellTracker Green, plated out on 384-Well plates (Greiner) at  $5 \times 10^3$  hMSC/cm<sup>2</sup> and cultured for 24 h. INA-6 cells were stained with CellTracker DeepRed, resuspended in MSC-medium, added to adhering hMSCs in different amounts ( $5 \times 10^3$ ,  $1 \times 10^3$ ,  $2 \times 10^3$  INA-6/cm<sup>2</sup>) and co-cultured for 24 h and 48 h. The complete co-culture was scanned and the number of INA-6 cells adhering on one MSC was counted manually for 100 MSCs for each technical replicate. Fluorescent images were digitally re-stained (INA-6 green, hMSC inverse black).

### **Analysis INA-6 Survival and Aggregation Depending on hMSCs Confluence**

To describe aggregate growth and survival of INA-6 depending on hMSC density, unstained hMSCs were seeded out into 96-well plates (white, clear bottom, Greiner) at different densities (Tab. S3). To ensure nutrient supply, we used lower cell densities for longer co-culturing durations while maintaining constant ratios of INA-6 to adhesion surface provided by hMSCs. Those plates that are to be assessed after 72 h of co-culturing received further 100 µL of fresh MSC-medium after 24 h of co-culturing (total volume of 300 µL), and after 48 h of co-culturing, 100 µL was removed gently from the co-culture and (carefully not to stir up co-culture on bottom) replaced with fresh MSC-Medium after 48 h of co-culturing.

To describe aggregate growth, complete wells were scanned using 10x magnification, phase contrast, 2x2 binning, and autofocus focusing on each tile both before and after harvesting. Afterwards, INA-6 cells were harvested for measuring viability and apoptosis.

**Tab. S3:** Seeding densities for describing growth and survival of INA-6 depending on hMSC density. Co-cult. dur. = Co culturing duration; MSC-adh. surface = adhesion surface provided by hMSCs; vol. = volume.

Co- cult. dur. [h]	hMSC density [1000 hMSC/cm <sup>2</sup> ]			INA-6 density [1000 INA-6/cm <sup>2</sup> ]	Ratios INA : MSC (adh. surface)			Seeding vol. [ $\mu$ L]	End vol. [ $\mu$ L]
24	2	10	40	10	1 : 0.2	1 : 1	1 : confluent	200	200
48	1	5	40	5	1 : 0.2	1 : 1	1 : confluent	200	200
72	1	5	40	5	1 : 0.2	1 : 1	1 : confluent	200 [after 24 h: + 100] [after 48 h: exchange 100]	300

For luminescent assessment of cell survival, INA-6 were harvested by removing co-culture medium, adding 150  $\mu$ L of MSC-Medium, and then stirred by strongly pipetting up and down twice while aiming the pipette tip at the upper corner, lower left and lower right of the well bottom ('Mercedes star'). Washing and stirring was repeated once before washing wells again with 150 mL MSC-Medium. Harvested INA-6 cells were strained (40  $\mu$ m), pelleted, and resuspended in 200  $\mu$ L MSC-Medium. Cells were counted using Neubauer chambers, re-distributed into 96-well plates (white, clear bottom) with  $1 \times 10^5$  INA-6 cells per well, and then subjected to viability and apoptosis assays.

To minimize the loss of sensitive apoptotic cells, another approach was used to measure viability and apoptosis without harvesting INA-6 cells. hMSCs and INA-6 were seeded out individually in parallel to the co-cultures (S. Tab. 02). Prior to measuring viability and apoptosis, culture volume was adjusted to 150  $\mu$ L by removing 50  $\mu$ L or 150  $\mu$ L for the timepoints 48 h or 72 h, respectively (carefully not to stir up culture on bottom). 100  $\mu$ L of luminescent reagents were then added directly to 150  $\mu$ L of co-culture. The fold change of viability or apoptosis that is due to MSC interaction ( $FC_{MSC\ interaction}$ ) was then calculated using the following formula, with  $L$  being the mean of four technical replicates measured in relative luminescent units per seconds [RLU/s],  $L_{Co\ culture}L_{MSC}$ ,  $L_{INA\ 6}$  the luminescence measured in the co-culture, hMSCs alone and INA 6 alone, respectively.

$$FC_{MSC\ Interaction} = \frac{L_{Co\ Culture}}{L_{MSC} + L_{INA\ 6}}$$

#### Time-Lapse Characterization of INA-6 Aggregation, Detachment and Division

In order to record the aggregation and detachment of INA-6 in contact with hMSCs, hMSCs (5e3 cells/cm<sup>2</sup>) were fluorescently stained with PKH26 and plated onto 8-well  $\mu$ -Slides (ibidi, Gräfelfing,

Germany). hMSCs were incubated for 24 h before being placed into an ibidi Stage Top Incubation System and were equilibrated to the incubation system for a minimum of 3 h (80% humidity and 5% CO<sub>2</sub>). INA-6 cells ( $2 \times 10^4$  cells/cm<sup>2</sup>) were washed and resuspended in 33% (v/v) MSC-conditioned medium before adding them directly before acquisition start in a small volume (10 µL). Brightfield and fluorescence images of 13 mm<sup>2</sup> of co-culture were acquired every 15 minutes for 63 h. Movement speed of the motorized table was adjusted to the lowest setting that allows acquisition of the complete region within 15 minutes.

Respective events of interest were analyzed manually and categorized into defined event parameters. Events were binned across the time axis using these boundaries: [0.0, 12.85, 25.7, 38.55, 51.4, 64.25]. We collected a minimum of events per recording and analysis so that each time bin contained at least 5 values, except when analyzing detachment events, since these did not appear before 20 h of incubation for some replicates. For each recording and event parameter, the event count was normalized by dividing by the total number of events per time bin.

We determined the frequency and the cause of aggregation by looking for two interacting INA-6 cells and went backward in time to see if they were two daughter cells or if two independent INA-6 cells had collided.

We determined the frequency of aggregates with detaching cells by tracing their growth across the complete time-lapse and looking for detachment events. We picked random 100 aggregates by including aggregates from both the border and center of the well.

We characterized detachment events by noting multiple parameters manually: Time point of detachment, aggregate size (at the time of detachment), the last interaction partner, and the number of detaching INA-6 cells.

For characterizing cell division events, we recorded a new set of time-lapse videos using unstained hMSCs that were grown to confluence for 24 h ( $4 \times 10^4$  hMSCs/cm<sup>2</sup>) to provide for unlimited adhesion surface. We categorized daughter cells in terms of their mobility (mobility being the speed of putative movements or “rolling”). The mobility criteria were met if one INA-6 daughter cell moved farther than half a cell radius within one frame (15 min) relative to the MSC-adherent INA-6 cell which was required to stand still in-between respective frames. We measured the “rolling” duration by subtracting the time point of the last perceived movement from the time point of division. We

excluded those division events from the measurement of rolling duration, if INA-6 cells underwent apoptosis shortly after division.

### **Cell Cycle Synchronization at M-Phase**

INA-6 cells were arrested at mitosis by double thymidine (2 mM) treatments followed by 5 h of nocodazole (500 ng/mL) incubation. In detail:  $3 \times 10^5$ /mL INA-6 in 4 mL were treated with 2 mM thymidine (Sigma) for 16.5 h. Cells were released by washing them in INA-6 medium once and allowed to cycle for 9 h before treating them with 2 mM thymidine for 18 h a second time. Afterwards, cells were released and allowed to cycle for 2 h before treating them with 100 ng/ mL nocodazole (Sigma) for 5 h. Arrested INA-6 were released by washing them once and resuspending them in MSC-medium with 33% MSC-conditioned medium. Cell cycle profile was checked using image cytometry (Fig S2).

### **V-Well Adhesion Assay**

This assay was modified from (Weetall et al., 2001). 96 v-well plates were coated with collagen I (rat tail, Corning). Collagen coating ensures that confluent hMSCs withstand centrifugation even after hMSCs in the well tip were removed. hMSCs ( $4 \times 10^4$  cells/cm<sup>2</sup>) were seeded out and grown to confluence for 24 h in collagen-coated v-well plates. To ensure that only INA-6 are pelleted in the v-well tip, hMSCs were removed from the well-tip by touching the well-ground with a 10 µL pipette and roughly pipetting hMSCs away.

Arrested INA-6 ( $1 \times 10^4$  cells/cm<sup>2</sup>) were released by washing them once in PBS and resuspending them in 33% (v/v) MSC-conditioned medium before adding them on top of confluent hMSCs. INA-6 adhered for 1, 2, 3 and 24 h before the complete co-culture was stained with 0.5 µM Calcein-AM (10 min at 37 °C).

Non-adherent INA-6 were pelleted by centrifugation using a Hettich 1460 rotor ( $r = 124$  mm) at 2000 rpm (555 g) for 10 min.

The well tip was imaged by fluorescence microscopy with 5x magnification, 96 HE emission filter, autofocus configured for maximum signal intensity, 2x2 binning and 14 bit grayscale depth. Pellet brightness was analyzed in ZEN 2.6 (Zeiss) by summing up pixel brightnesses across the complete pellet image. Background brightness was acquired from a cell culture with only hMSCs. Reference brightness was acquired from a cell culture with only INA-6, defining 100% pellet brightness without

adhesion. Background intensity was subtracted before normalizing by reference. Outliers were removed from technical replicates ( $n=4$ ) if their z-score was larger than  $1.5 \sigma$  technical variation.

After measuring pellet brightnesses, the cell pellet was removed by pipetting 10  $\mu\text{L}$  from the well tip.

Pellets of the same technical replicates were pooled, washed in PBS, resuspended in 200  $\mu\text{L}$  PBS, added to 1.8 mL ice-cold 70% ethanol and stored at -20 °C.

Remaining non-MSC-adhering INA-6 cells were removed by replacing culture medium with 100  $\mu\text{L}$  of medium. MSC-adherent INA-6 were manually detached by rapid pipetting and equally pelleted, analyzed, and isolated.

### **Cell Cycle Profiling**

INA-6 cells were fixed in 70% ice-cold ethanol, washed, resuspended in PBS, distributed in 96-well plates and stained with Hoechst-33342 (2.5  $\mu\text{g/mL}$  in PBS) for 1 h at 37 °C.

For image cytometric cell cycle profiling, plates were scanned completely using automated fluorescence microscopy with 5x magnification, 96 HE emission filter, 1x1 binning, 14 bit depth and an illumination time that fills 70% of grayscale range. The autofocus was configured to re-adjust every second tile. A pre-trained convolutional neural network (“DeepFeatures 2 reduced”, Intellesis, Zeiss) was fine-tuned to segment scans into background, single nuclei and fragmented nuclei. Nuclei were filtered to exclude fragmented nuclei and those nuclei with extreme size (within the range of 50-500  $\mu\text{m}^2$ ) and roundness (within the range of 0.4-1.0). Cell cycle profiles were normalized by the mode of the nucleus intensities within the G0/G1 peak. To retrieve frequencies of cells cycling in G0/G1, S, and G2 phase, the brightness distribution of all single nuclei was fitted to the sum of three Gaussian curves (“Skewed Gaussian Model” for G0G1 and G2 phase, and “Rectangle Model” for S phase) using the python package LMFIT (Newville et al., 2014) (Fig. S4). The gaussian curves were used to calculate the cell frequencies for each cell cycle phase by integration using the composite trapezoidal rule implemented by numpy.trapz (Harris et al., 2020).

For validation of image cytometry, 5 mL of INA-6 stock culture was removed and ethanol fixed as described above. Flow cytometry analyses were performed using an Attune Nxt Flow Cytometer (Thermo Fisher, USA). Data analyses were performed using FlowJo V10 software (TreeStar, USA).

**Protocol: Well Plate Sandwich Centrifugation (WPSC)**

96 well plates (flat bottom, clear) were coated with collagen I (rat tail, Corning). Collagen coating ensures that confluent hMSCs withstand centrifugation and repeated washing. hMSCs ( $2 \times 10^4$  cells/cm<sup>2</sup>) were seeded out and grown to confluence for 72 h in collagen-coated 96-well plates.

To remove aggregates from the medium and prevent clogging of magnetic columns, we strained any FCS-containing fluid with a 40 µm cell strainer.

Collect MSC-conditioned medium and add INA-6:

1. Collect hMSC-conditioned medium (CM) from the well plates and replace it with 100 µL of fresh hMSC medium. Collect CM from different donors separately
2. Strain CM (40 µm) and centrifuge it for 15 minutes at 250 g to ensure that CM does not contain hMSCs
3. Dilute CM by mixing 2 parts of CM with 1 part of MSC-medium (dilute 1.5 fold)
4. Count INA-6 cells and retrieve enough cells to fill all 96 wells with  $2 \times 10^4$  INA-6/cm<sup>2</sup> ( $6.8 \times 10^4$  cells per well, covering ~65% of the well bottom).
5. Centrifuge INA-6 (5 min, 250 g) and resuspend them in a volume of diluted CM to reach a concentration of  $6.8 \times 10^5$  INA-6/mL
6. Add 100 µL INA-6 suspension to hMSCs (end volume: 200 µL; end concentration: 33% (v/v) hMSC-conditioned medium)
7. Incubate for 24 h at 37 °C and 5% CO<sub>2</sub>

Prepare CM-INA6 reference:

8. Add 100 µL of fresh MSC-medium into each well of an empty 96-well plate (not coated)
9. Add 100 µL of INA-6 suspension ( $6.8 \times 10^5$  INA-6/mL in diluted CM)
10. Incubate for 24 h at 37 °C and 5% CO<sub>2</sub>

Collect CM-INA6 and nMA-INA6

11. Pre-warm well plate centrifuge to 37 °C
12. Prepare a counter-weight by filling 200 µL of water into all wells of an empty 96-well plate
13. Prepare well-plate sandwiches:

- a. Turn an empty 96-well plate (“catching plate”) upside down and place one on top of the co-culture-plate, the CM-IN6 reference plate, and the counter-weight so that all well openings align.
  - b. Fix well plates using tape with reusable adhesive (e.g. Leukofix)
14. Turn both plates around. Medium will spill from the co-culture plate into the catching plate
  15. Centrifuge plate for 40 seconds at 1000 rpm with the catching plate facing the ground
  16. Remove the adhesive tape and the co-culture plate.
  17. Turn the co-culture plate around and add 30 µL of washing medium (MSC-Medium 0% FCS, 3 mM EDTA) gently by touching the wall of each well and pressing the pipette slowly.
    - a. *Work quickly to ensure that co-culture does not dry. We recommend using a multipette (Eppendorf).*
    - b. *Many nMA-IN6 are removed by physical force applied by adding 30 µL of medium and not just by centrifugation. Hence, it is critical to apply the same dispensing technique across all replicates. We recommend using a multipette (Eppendorf) that can apply 30 µL with controllable pressure, since its push-button retains a long pushing path even for dispensing small volumes, unlike push-buttons from the usual 100 µL pipettes that reduce the pushing-path for smaller volumes.*
    - c. *Centrifugation minimizes technical variability by replacing one step of manual pipetting. Also, it ensures that confluent MSCs remain unharmed. Manual pipetting on the other hand would require touching the well-bottom to remove all fluids which damages the adhesive hMSC layer.*
  18. Turn the co-culture plate upside down, place it onto the catching plate and re-apply adhesive tape to fix the wellplate sandwich
  19. Repeat steps 14-18 two more times until the catching plate contains 290 µL of medium in each well
  20. Pool CM-IN6 from the catching plate that was fixed to the reference plate
  21. Pool nMA-IN6 from the catching plate that was fixed to the co-culture plate
  22. Collect remaining IN6 by adding 100 µL of PBS into each well of the catching plates, collect and pool with CM-IN6 or nMA-IN6.

23. Strain CM-INA6 and nMA-INA6 using 40 µm cell strainer
  24. Isolate MA-INA6 by continue with either accutase dissociation or rough pipetting
- Collect MA-INA6 by accutase dissociation followed by MAC sorting
25. Block 2 mL tubes with sorting buffer (PBS, 2 mM EDTA, 1% BSA) for 1 h at 4 °C
  26. Dilute accutase (Sigma A6964) (400-600 units/mL) 4-fold in cold PBS. Always keep accutase on ice, since accutase loses activity at room temperature.
  27. Add 50 µL of cold accutase (directly after the last centrifugation step) and incubate co-culture plate for 5 minutes at 37 °C.
  28. Place a co-culture plate onto a shaker and shake for 1 minute at 300 rpm.
  29. Collect cell suspension from wells and stop the reaction by adding 500 µL of FCS to pooled cell suspension.
  30. Evaluate presence of adherent INA-6 cells and the integrity of confluent hMSCs under the microscope.
  31. Repeat steps 24-27 until all INA-6 cells have dissociated or until confluent hMSCs start to tear.
  32. Strain cell suspension (30 µm). This yields MA-MSC.
  33. Pellet MA-INA6, nMA-INA6 and CM-INA6 (1200 rpm, 10 min).
  34. Resuspend MA-INA6 in 86 µL sorting buffer (PBS, 2 mM EDTA, 1% BSA)
  35. Resuspend CM-INA6 and nMA-INA6 in 300 µL cold diluted accutase and incubate for 3 min at 37 °C to ensure equal treatment for all samples.
  36. Stop accutase by adding 200 µL of FCS (100%)
  37. Pellet CM-INA6 and nMA-INA6 (1200 rpm, 10 min) and resuspend in 86 µL sorting buffer (PBS, 2 mM EDTA, 1% BSA).
  38. Transfer samples into 2 mL tubes that were blocked with sorting buffer
  39. Add 10 µL of CD45 coated magnetic beads (Miltenyi Biotec B.V. & Co. KG, Bergisch Gladbach)
  40. Place tubes into rotator and incubate for 15 minutes at 4 °C
  41. Continue with MAC sorting according to the manual. Use an MS column and wash 3 times.

42. Improve purity of eluted MA-INA6 by straining eluate ( $30\ \mu\text{m}$ ) (wash strainer using 1 mL of sorting buffer) and applying it onto an MS column a second time. Wash three times.

43. Collect 20  $\mu\text{L}$  per eluate and apply it onto a 96-well plate to evaluate purity

a. Incubate plate for 24 h

b. Count the number of adherent cells (hMSCs) per INA-6 using phase contrast microscopy

c. *We reached a mean purity of  $3.2 \times 10^{-4}$  ( $\pm 2.2 \times 10^{-4}$ ) hMSCs per MA-INA6.*

d. *hMSC contamination did not have an impact on RNAseq, since those genes that are highly expressed in hMSCs (VCAM1, ALPL, FGF5, FGFR2), did not appear as differentially expressed in MA-INA6 (Data not shown). RNAseq detected  $0.44 \pm 0.16$  CPM-normalized counts of VCAM1 transcripts in MA-INA6, however, it was excluded like all genes with less than 1 count in at least 2 of 5 replicates.*

44. Count cells using a Neubauer chamber

45. Pellet samples (250 g for 5 min)

46. Resuspend in respective medium or lysis buffer (e.g. RA1 for RNA extraction)

Collect MA-INA6 by rough pipetting (no MAC sorting)

47. After the last centrifugation step, add hMSC-medium to each well of the co-culture plate to reach a volume of 150  $\mu\text{L}$

a. *Since the yield of MA-INA6 was large, we dissociated MA-INA-6 cells from hMSCs by vigorous pipetting (for further samples after RNAseq, see Tab. S1). Since no enzymatic digestion is used, we reckoned that there would be no need for MAC sorting. Confluent hMSCs withstand this procedure and don't dissociate as single cells, which can be removed by straining cells ( $30\ \mu\text{m}$ ). We reached similar purities as for MAC-sorting (Data not shown).*

48. Using a multi-channel pipette (100  $\mu\text{L}$ ), gently raise 90  $\mu\text{L}$  into the tips

49. Lean pipette tip on the upper well-border and roughly pipette up and down once

50. Repeat step 48 at the lower right and lower left well border (Total of 3 pipetting steps  
“Mercedes Star”)

51. Attach a catching plate onto the co-culture and centrifuge for 40 seconds at 500 rpm (28 g)

52. Repeat steps 46-50 until a sufficient amount of MA-INA6 is removed
53. Control purity of MA-INA-6 by placing out aliquot onto an empty 96-well plate.
54. Collect MA-INA6 from catching plate
55. Remove hMSCs by straining cell suspension (30 µm)
56. Count cells using a Neubauer chamber
57. Pellet MA-INA6 (250 g for 5 min)
58. Resuspend in respective medium or lysis buffer

Centrifugal force: We used a Hettich 1460 rotor ( $r = 124$  mm) (Hettich GmbH & Co. KG, Tuttlingen, Germany). For calculating the centrifugal force that acts onto the co-culture within well plate sandwiches, we subtracted the height of the catching plate (14.4 mm, Greiner 96 well plate) and the depth of each well (10.9 mm). This yields a radius of 98.7mm, which translates to the following centrifugal forces: 500 rpm: 28 g; 1000 rpm: 110 g; 2000 rpm: 441 g.

Washing medium with EDTA: EDTA removes calcium from integrins which are required for adhesion. It is not strong enough to dissociate INA-6 from hMSCs, but could help with removing INA-6 from other INA-6. For generating samples for RNAseq, we added 3 mM of EDTA to washing medium. For further samples, we did not add EDTA to the washing medium, since we found that it does not increase yield for all biological replicates consistently (Data not shown). We suspect that integrin-mediated adhesion depends on hMSC donor or internal variance of INA-6. We recommend using 3 mM of EDTA, however, this requires further optimizations like including an incubation time at 37 °C after the addition of washing medium to account for biological variance. However, this could take long incubation times of up to 60 minutes (Lai et al., 2022).

### Track Cell Number During WPSC

To track the cell count during WPSC, INA-6 were stained with CellTracker green and both in co-culturing- and catching plates were scanned after each centrifugation step. For each round of centrifugation, an empty catching plate was used. A pre-trained convolutional neural network (Intellesis, Zeiss) was fine-tuned to segment the scans into background, cells, and cell borders. Single cells were counted and the cumulative sum for each catching plate was calculated.

### **Sub-Culturing After WPSC of MSC-Interacting INA-6 Subpopulations**

After CM-INA6, nMA-INA6, and MA-INA6 were isolated, they were counted with a Neubauer chamber using all nine quadrants and diluted to  $10^5$  cells/mL in MSC-medium (10% FCS, no IL-6 except for control). 100 µL of cell suspension was applied to 96-well plates, incubated for 48 h at 37 °C and 5% CO<sub>2</sub> and then subjected to viability and apoptosis assays.

### **RNA Isolation**

Total RNA was isolated from INA-6 cells by using the NucleoSpin RNA II Purification Kit (Macherey-Nagel, Düren, Germany) according to the manufacturer's instructions.

### **RNAseq, Differential Expression and Functional Enrichment Analysis of INA-6 cells**

FASTQ files were merged to the respective sample. The quality of FASTQ files was assessed with FastQC (Andrews, 2010) tool, and a joint report was created with MultiQC (Ewels et al., 2016) tool. Fastq files were aligned with STAR (Dobin et al., 2013) to the GRCh38 reference genome build (Zerbino et al., 2018). Quality and alignment statistics of final BAM files were assessed with samtools stats (Li et al., 2009), and a joint report with FastQC reports by MultiQC was generated.

Raw read counts were generated with HTSeq (Anders et al., 2015) with the union method. HTSeq runs internally in STAR. Differential gene expression analysis was done with edgeR (Robinson et al., 2010) in R 3.6.3 (R Core Team, 2018), according to the edgeR manual.

Counts were merged and genes with zero counts in all samples were removed (number of genes: 36380).

The whole count table was annotated with R Bioconductor (Gentleman, n.d.) (Gentleman et al. 2004) human annotation data package org.Hs.eg.db (Carlson, 2016).

A DGEList Element was created with the raw counts, gene information, i.e. Ensembl GenelIDs, HUGO Symbol, Genename, and ENTREZ GenelIDs and a sample grouping meta data table.

```
y <- DGEList(counts=ct2[,-1:4], group=meta.data$group, genes=ct2[,1:4])
```

Counts were filtered to keep only those genes which have at least 1 read per million in at least 2 samples (number of genes: 14136). Afterwards normalization factors were recalculated.

```
keep <- rowSums(cpm(y)>1) >=2  
y <- y[keep, , keep.lib.size = FALSE]  
y1 <- calcNormFactors(y)
```

A design matrix was created with grouping factor by treatment condition (group=F1, F2, F3, which are abbreviations for CM-INA6, nMA-INA6, MA-INA6, respectively)

```
design = model.matrix(~0+group)
```

Dispersion was estimated, the resulting coefficient of biological variation (BCV) is 0.135, i.e. BCV expression values vary up and down by 13.5% between samples.

```
y1.1 <- estimateDisp(y1, design)  
BCV <- sqrt(model.F$y1.1$common.dispersion)
```

A generalized linear (glmQLFit function) model was fitted.

```
fit <- glmQLFit(y1.1, design)
```

and pairwise comparisons were made, e.g.

```
F1vsF2 <- glmQLFTest(fit, contrast = makeContrasts(groupF1 - groupF2,  
levels = design))
```

top significant differential expressed genes were written to a table

```
DE.F1vsF2 <- topTags(F1vsF2, n=nrow(F1vsF2), p.value = 0.05)
```

Afterwards, gene list of differentially expressed genes were used for functional enrichment analysis with metascape (Zhou et al., 2019).

### RT-qPCR

For cDNA synthesis 1 µg of total RNA was reverse transcribed with Oligo(dT)15 primers and Random Primers (both Promega GmbH, Mannheim, Germany) and Superscript IV reverse transcriptase (Thermo Fisher Scientific) according to the manufacturer's instructions. For quantitative PCR the cDNA was diluted 1:10 and qPCR was performed in 20 µl by using 2 µl of cDNA and 10 µl of GoTaq qPCR Master Mix (Promega GmbH) and 5 pmol of sequence-specific primers obtained from biomers.net GmbH (Ulm, Germany) or Qiagen GmbH (Hilden, Germany) (see Tab. S4 for primer sequences and PCR conditions). qPCR conditions were as follows: 95°C for 3 min; 40 cycles: 95°C for 10 s; respective annealing temperature for 10 s; 72°C for 10 s; followed by melting curve analysis for the specificity of qPCR products by using the qPCR thermal cycler Professional Thermocycler Biometra (Analytik Jena AG, Jena, Germany). Samples that showed unspecific byproducts were discarded. Ct values were measured in three technical replicates (triplicates). Non-detects were discarded. One of three technical replicates was treated as an outlier and excluded if

its z-score crossed  $1.5\sigma$  technical variation. We normalized expression by the housekeeping gene 36B4. Efficiencies were determined in each reaction by linear regression of log transformed amplification curve (Ramakers et al., 2003). Differential expression was calculated based on a modified  $\Delta\Delta Ct$  formula that separated exponents to apply individual efficiencies to each Ct value:

$$\text{Fold Change} = \frac{E_{tar}^{\Delta Ct_{tar}(co-treated)}}{E_{ref}^{\Delta Ct_{ref}(co-treated)}} = \frac{E_{tar,co}^{Ct_{tar,co}} : E_{tar,treated}^{Ct_{tar,treated}}}{E_{ref,co}^{Ct_{ref,co}} : E_{ref,treated}^{Ct_{ref,treated}}}$$

$E_{tar,co}$  = Efficiency of the target gene measured in the control sample

$Ct_{tar,co}$  = Ct value of the target gene measured in the control sample

$tar$  = Target Gene;  $ref$  = Reference Gene

$treated$  = Treated sample;  $co$  = Control Sample

Fold change expression was normalized by the median of CM-INA6 (and not samplewise, as commonly used in  $\Delta\Delta Ct$ ) since some genes were not expressed without direct MSC contact s (e.g. MMP2), and also in order to display variation of CM-INA6 next to nMA-INA6 and MA-INA6.

**Tab. S4:** List of primers. Some primers required a melting step to be performed before fluorescent readout to remove byproducts.

Primer	Sequence 5' - 3'	base pairs [bp]	annealing temp. [°C]
36B4_s	tgcatcagtacccatttatcat	122	60
36B4_as	aggcagatggatcagccaaga		
BCL6_s	tagagcccataaaacggtcctcat	221	55 + Melting Step at 77 °C
BCL6_as	cgc当地点attgagccgagatgtgt		
BMP4_s	tacatgcgggatcttaccg	132	58
BMP4_as	atgttcttcgtggtaagc		
BTG2_s	gtattcttgtagggccacactaa	264	60 + Melting Step at 78 °C
BTG2_as	tcttaagggtgattcggtttggaa		
CXCL8_s	actgagagtgattgagagtggacc	251	55 + Melting Step at 77 °C
CXCL8_as	ccctacaacagacccacacaatac		
CXCL12_s	gattcttcgaaagccatgttgcga	119	56

CXCL12_as	caatgcacacttgtctgttgtgt		
DCN_s	caacaacaagcttaccagagtacct	160	57
DCN_as	tgaaaagactcacacccgaataaga		
DKK1_s	gcactgatgagtagtactgcgctag	129	56
DKK1_as	ttttgcagtaattccggggc		
IL10RB_s	gagtgagcctgtctgtgagcaa	139	55
IL10RB_as	cttgtaaacgcaccacagcaag		
IL24_s	caaacagttggacgtagaaggcagc	149	55
IL24_as	tgaaatgacacagggAACAAACCA		
LOX_s	ctgctcagattccccaaag	125	57
LOX_as	tggcatcaagcaggtaatcgat		
MMP2_s	ttgtatttgatggcatcgctcaga	155	56
MMP2_as	cgtataccgcatcaatctttccg		
MMP14_s	cgacaagattgtatgcgtc	140	57
MMP14_as	tcccttcccagactttgatg		
MUC1_s	gcagcctctcgatataacctg	200	58
MUC1_as	gtaggtgggtactcgctca		
NOTCH2_s	gtgcttgttgaacacttgtgcc	185	55
NOTCH2_as	cactcgcatctgtatccaccaatg		
OPG (TNFRSF11B)	no sequence available (Proprietary primers from Qiagen: QT00014294 TNFRSF11B_1_SG)		60
PRICKLE1_s	cagaggtatatcatgaaggacggc	102	56
PRICKLE1_as	gtcccacaccaatatgttccccac		
TGM2_s	caaccttctcatcgagacttccg	100	58
TGM2_as	tcatccacgactccacccag		
TNFRSF1A_s	ctccttcaccgcttcagaaaacc	153	55
TNFRSF1A_as	ttcactccaataatgccggtaactg		
TRAF5_s	tgcctgttagataaagaggcatca	177	56
TRAF5_as	aacactgcacaggttcaaataagc		

**Statistics**

For molecular analyses, each data point represents one biological replicate, which we define as the mean of all technical replicates of co-cultures that were seeded out from the same batch of hMSCs and/or INA-6 cells on the same day. For analyses of time-lapse recordings, each datapoint represents the normalized event count from a recording of one co-culture. We prioritized unique hMSCs for each biological replicate or recording (Tab. S1). Bars and lines represent the mean and error bars represent the standard deviation of all hMSC donors or recordings (= all biological replicates).

Metric, normal distributed, dependent data was analyzed using factorial RM-ANOVA and paired Student's t-test. Results of RM-ANOVA are reported as such:  $[F(df_1, df_2) = F; p = p\text{-value}]$ , with  $df_1$  being the degrees of freedom of the observed effect,  $df_2$  being the degrees of freedom of the error and  $F$  being the F-statistic (Vallat, 2018). If sphericity was met, p-values were not corrected with the Greenhouse-Geisser method (p-unc).

$$df_1 = k - 1 \quad k = \text{The number of groups (of a factor, if factorial RM-ANOVA)}$$

$$df_2 = (k - 1)(n - 1)n = \text{The number of samples in each group}$$

$$F = \frac{SS_{Effect} \div df_1}{SS_{Error} \div df_2} \quad SS = \text{Sums of squares for effect or error}$$

If datapoints within dependent sample pairs were missing, such pairs were excluded from paired t-test while other pairs of the same subject remained.

Metric non-normal distributed, independent data was analyzed using Kruskal-Wallis H-test and Mann-Whitney U tests. Results of Kruskal-Wallis H-test was reported as such:  $[H(df) = H]$ , with  $df$  being the degrees of freedom and  $H$  being the Kruskal-Wallis H statistic, corrected for ties (Vallat, 2018).

$$df = k - 1 \quad k = \text{The number of groups}$$

Metric bivariate non-normal distributed data was correlated using spearman's rank correlation and reported as such:  $[\rho(df) = \rho, p = p\text{-value}]$ , with  $\rho$  being Spearman's rank correlation coefficient.  $df$  is calculated as such:

$$df = n - 2 \quad n = \text{The number of observations}$$

These test were applied using the python (3.10) -packages pingouin (0.5.1). For three-factor RM-ANOVA we used statsmodels (0.14.0) (Seabold & Perktold, 2010; Vallat, 2018). Data was plotted

using seaborn (Waskom, 2021). Sphericity was ensured by Mauchly's test. Normality was checked with the Shapiro-Wilk test for  $n > 3$ .

Datapoints were log10 transformed to convert the scale from multiplicative ("foldchange") to additive, or in order to fulfill sphericity requirements.

P-values derived from patient survival data were corrected using the Benjamini-Hochberg procedure. For other post-hoc analyses, p-values were not adjusted for family-wise error rate in order to minimize type I errors. To prevent type II errors, the same conclusions were validated by different experimental setups and through varying hMSCs donors across experiments (Tab. S1).

Significant p-values from pairwise tests were annotated as stars between data groups (p-value: 0.05 > \* > 0.01 > \*\* >  $10^{-3}$  > \*\*\*  $10^{-4}$  > \*\*\*\*). If too many significant pairs were detected, we annotated only those pairs of interest.

No power calculation was performed to determine sample size since samples were limited by availability of primary hMSC donors. Experiments were repeated until a minimum of three biological replicates were gathered.

### **Patient Cohort, Analysis of Survival and Expression**

Patient samples ( $n=873$ ) were collected at the UKHD and processed as described (Seckinger et al., 2017, 2018), and are available at the European Nucleotide Archive (ENA) via accession numbers PRJEB36223 and PRJEB37100. Consecutive patients with monoclonal gammopathy of unknown significance (MGUS) ( $n = 62$ ), asymptomatic ( $n = 259$ ), symptomatic, therapy-requiring ( $n = 764$ ), and relapsed/refractory myeloma ( $n = 90$ ), as well as healthy donors ( $n = 19$ ) as comparators were included in the study approved by the ethics committee (#229/2003, #S-152/2010) after written informed consent.

Gene expression was measured by RNA sequencing as previously described (Seckinger et al., 2018). Gene expression is defined as the log2 transformed value of normalized counts + 1 (as pseudocount). Progression-free (PFS) and overall survival (OS) was analyzed for the subset of previously untreated symptomatic MM patients. For delineating "high" and "low" expression of target adhesion ( $n=101$ ) and cell cycle ( $n=173$ ) genes, thresholds per gene were calculated with maximally selected rank statistics by the maxstat package in R (Hothorn & Lausen, n.d.). PFS and OS were

analyzed for high vs. low expression with the Kaplan-Meier method (Kaplan & Meier, 1958). Significant differences between the curves were analyzed with log-rank tests (Harrington & Fleming, 1982). P-values were corrected for multiple testing by the Benjamini-Hochberg method. Analyses were performed with R version 3.6.3 (R Core Team, 2018).

## References

- Anders, S., Pyl, P. T., & Huber, W. (2015). HTSeq—A Python framework to work with high-throughput sequencing data. *Bioinformatics (Oxford, England)*, 31(2), 166–169.  
<https://doi.org/10.1093/bioinformatics/btu638>
- Andrews, S. (2010). FASTQC. A quality control tool for high throughput sequence data.
- Burger, R., Guenther, A., Bakker, F., Schmalzing, M., Bernand, S., Baum, W., Duerr, B., Hocke, G. M., Steininger, H., Gebhart, E., & Gramatzki, M. (2001). Gp130 and ras mediated signaling in human plasma cell line INA-6: A cytokine-regulated tumor model for plasmacytoma. *The Hematology Journal: The Official Journal of the European Haematology Association*, 2(1), 42–53. <https://doi.org/10.1038/sj.thj.6200075>
- Carlson, M. (2016). Org.Hs.eg.db. *Bioconductor*. <https://doi.org/10.18129/B9.bioc.Org.Hs.eg.db>
- Chatterjee, M., Hönenmann, D., Lentzsch, S., Bommert, K., Sers, C., Herrmann, P., Mathas, S., Dörken, B., & Bargou, R. C. (2002). In the presence of bone marrow stromal cells human multiple myeloma cells become independent of the IL-6/gp130/STAT3 pathway. *Blood*, 100(9), 3311–3318. <https://doi.org/10.1182/blood-2002-01-0102>
- Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., & Gingeras, T. R. (2013). STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics*, 29(1), 15–21. <https://doi.org/10.1093/bioinformatics/bts635>
- Ewels, P., Magnusson, M., Lundin, S., & Käller, M. (2016). MultiQC: Summarize analysis results for multiple tools and samples in a single report. *Bioinformatics*, 32(19), 3047–3048.  
<https://doi.org/10.1093/bioinformatics/btw354>
- Fernandez-Rebollo, E., Mentrup, B., Ebert, R., Franzen, J., Abagnale, G., Sieben, T., Ostrowska, A., Hoffmann, P., Roux, P.-F., Rath, B., Goodhardt, M., Lemaitre, J.-M., Bischof, O., Jakob, F., & Wagner, W. (2017). Human Platelet Lysate versus Fetal Calf Serum: These Supplements Do Not Select for Different Mesenchymal Stromal Cells. *Scientific Reports*, 7, 5132. <https://doi.org/10.1038/s41598-017-05207-1>
- Gentleman. (n.d.). *Bioconductor—BiocViews*. Retrieved June 9, 2023, from <https://bioconductor.org/packages/3.17/BiocViews.html>

- Gramatzki, M., Burger, R., Trautman, U., Marschalek, R., Lorenz, H., Hansen-Hagge, T. E., Baum, W., Bartram, C. R., Gebhart, E., & Kalden, J. R. (1994). *Two new interleukin-6 dependent plasma cell lines carrying a chromosomal abnormality involving the IL-6 gene locus.* 84 Suppl. 1, 173a–173a.
- Greenstein, S., Krett, N. L., Kurosawa, Y., Ma, C., Chauhan, D., Hidemitsu, T., Anderson, K. C., & Rosen, S. T. (2003). Characterization of the MM.1 human multiple myeloma (MM) cell lines: A model system to elucidate the characteristics, behavior, and signaling of steroid-sensitive and -resistant MM cells. *Experimental Hematology*, 31(4), 271–282.  
[https://doi.org/10.1016/s0301-472x\(03\)00023-7](https://doi.org/10.1016/s0301-472x(03)00023-7)
- Harrington, D. P., & Fleming, T. R. (1982). A Class of Rank Test Procedures for Censored Survival Data. *Biometrika*, 69(3), 553–566. <https://doi.org/10.2307/2335991>
- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., ... Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585(7825), Article 7825.  
<https://doi.org/10.1038/s41586-020-2649-2>
- Hothorn, T., & Lausen, B. (n.d.). *Maximally Selected Rank Statistics in R* [Computer software].  
<http://cran.r-project.org/web/packages/maxstat/index.html>.
- Kaplan, E. L., & Meier, P. (1958). Nonparametric Estimation from Incomplete Observations. *Journal of the American Statistical Association*, 53(282), 457–481.  
<https://doi.org/10.1080/01621459.1958.10501452>
- Lai, T.-Y., Cao, J., Ou-Yang, P., Tsai, C.-Y., Lin, C.-W., Chen, C.-C., Tsai, M.-K., & Lee, C.-Y. (2022). Different methods of detaching adherent cells and their effects on the cell surface expression of Fas receptor and Fas ligand. *Scientific Reports*, 12(1), Article 1.  
<https://doi.org/10.1038/s41598-022-09605-y>
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., & 1000 Genome Project Data Processing Subgroup. (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25(16), 2078–2079.  
<https://doi.org/10.1093/bioinformatics/btp352>

- Newville, M., Stensitzki, T., Allen, D. B., & Ingargiola, A. (2014). *LMFIT: Non-Linear Least-Square Minimization and Curve-Fitting for Python* [Computer software]. Zenodo.  
<https://doi.org/10.5281/zenodo.11813>
- Nilsson, K., Bennich, H., Johansson, S. G., & Pontén, J. (1970). Established immunoglobulin producing myeloma (IgE) and lymphoblastoid (IgG) cell lines from an IgE myeloma patient. *Clinical and Experimental Immunology*, 7(4), 477–489.
- R Core Team. (2018). *R: A language and environment for statistical computing* [Manual]. R Foundation for Statistical Computing. <https://www.R-project.org/>
- Ramakers, C., Ruijter, J. M., Deprez, R. H. L., & Moorman, A. F. M. (2003). Assumption-free analysis of quantitative real-time polymerase chain reaction (PCR) data. *Neuroscience Letters*, 339(1), 62–66. [https://doi.org/10.1016/S0304-3940\(02\)01423-4](https://doi.org/10.1016/S0304-3940(02)01423-4)
- Robinson, M. D., McCarthy, D. J., & Smyth, G. K. (2010). EdgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics (Oxford, England)*, 26(1), 139–140. <https://doi.org/10.1093/bioinformatics/btp616>
- Seckinger, A., Delgado, J. A., Moser, S., Moreno, L., Neuber, B., Grab, A., Lipp, S., Merino, J., Prosper, F., Emde, M., Delon, C., Latzko, M., Gianotti, R., Lüoend, R., Murr, R., Hosse, R., J., Harnisch, L. J., Bacac, M., Fauti, T., ... Vu, M. D. (2017). Target Expression, Generation, Preclinical Activity, and Pharmacokinetics of the BCMA-T Cell Bispecific Antibody EM801 for Multiple Myeloma Treatment. *Cancer Cell*, 31(3), 396–410. <https://doi.org/10.1016/j.ccr.2017.02.002>
- Seckinger, A., Hillengass, J., Emde, M., Beck, S., Kimmich, C., Dittrich, T., Hundemer, M., Jauch, A., Hegenbart, U., Raab, M.-S., Ho, A. D., Schönland, S., & Hose, D. (2018). CD38 as Immunotherapeutic Target in Light Chain Amyloidosis and Multiple Myeloma-Association With Molecular Entities, Risk, Survival, and Mechanisms of Upfront Resistance. *Frontiers in Immunology*, 9, 1676. <https://doi.org/10.3389/fimmu.2018.01676>
- Vallat, R. (2018). Pingouin: Statistics in Python. *Journal of Open Source Software*, 3(31), 1026. <https://doi.org/10.21105/joss.01026>
- Weetall, M., Hugo, R., Maida, S., West, S., Wattanasin, S., Bouhel, R., Weitz-Schmidt, G., Lake, P., & Friedman, C. (2001). A Homogeneous Fluorometric Assay for Measuring Cell

Adhesion to Immobilized Ligand Using V-Well Microtiter Plates. *Analytical Biochemistry*, 293(2), 277–287. <https://doi.org/10.1006/abio.2001.5140>

Zerbino, D. R., Achuthan, P., Akanni, W., Amode, M. R., Barrell, D., Bhai, J., Billis, K., Cummins, C., Gall, A., Girón, C. G., Gil, L., Gordon, L., Haggerty, L., Haskell, E., Hourlier, T., Izuogu, O. G., Janacek, S. H., Juettemann, T., To, J. K., ... Flórek, P. (2018). Ensembl 2018. *Nucleic Acids Research*, 46(D1), D754–D761. <https://doi.org/10.1093/nar/gkx1098>

Zhou, Y., Zhou, B., Pache, L., Chang, M., Khodabakhshi, A. H., Tanaseichuk, O., Benner, C., & Chanda, S. K. (2019). Metascape provides a biologist-oriented resource for the analysis of systems-level datasets. *Nature Communications*, 10(1), Article 1. <https://doi.org/10.1038/s41467-019-09234-6>

# Chapter 2: Semi-Automation of Data Analysis

## Abstract

plotastic addresses the challenges of transitioning from exploratory data analysis to hypothesis testing in Python’s data science ecosystem. Bridging the gap between seaborn and pingouin, this library offers a unified environment for plotting and statistical analysis. It simplifies the workflow with user-friendly syntax and seamless integration with familiar seaborn parameters (`y`, `x`, `hue`, `row`, `col`). Inspired by seaborn’s consistency, plotastic utilizes a `DataAnalysis` object to intelligently pass parameters to pingouin statistical functions. Hence, statistics and plotting are performed on the same set of parameters, so that the strength of seaborn in visualizing multidimensional data is extended onto statistical analysis. In essence, plotastic translates seaborn parameters into statistical terms, configures statistical protocols based on intuitive plotting syntax and returns a `matplotlib` figure with known customization options and more. This approach streamlines data analysis, allowing researchers to focus on correct statistical testing and less about specific syntax and implementations.

## Introduction

The reproducibility crisis in research highlights a significant challenge in contemporary bio-sciences, where a substantial portion of studies faces reproducibility issues (Begley & Ioannidis, 2015). One critical yet often overlooked aspect contributing to this crisis is data management. The literature most often refers to *big data* as the main challenge (Gomez-Cabrero et al., 2014). However, these challenges are also present in smaller datasets, which the author refers to as *semi-big data*. This term describes datasets that, while not extensive enough to necessitate advanced computational tools typically reserved for *big data*, are sufficiently large to render manual analysis very time-intensive. *Semi-big data* is often generated by methods like automated microscopy or multiplex qPCR, which produce volumes of data that are manageable on a surface level, but pose substantial barriers for in-depth, manual reproducibility (Bustin, 2014). This is further complicated by the complexity inherent in multidimensional datasets. For example, the qPCR experiment from Chapter 1, Fig. 4 involves the analysis of 19 genes across three subpopulations, including eleven biological and three technical replicates, resulting in a total of 1881 data points that are all assigned to a complex set of experimental variables. Without a clearly documented data analysis protocol and standardized data formats, the reproduction of such analysis becomes extremely challenging, if not impossible (Bustin, 2014).

The evolving standards in data analysis advocate for the standardization of analytical pipelines, rationalization of sample sizes, and enhanced infrastructure for data storage, address-

ing some of these challenges (Goodman et al., 2016; Wilkinson et al., 2016). However, these advancements can place undue pressure on researchers, particularly those with limited training in statistics, underscoring the need for intuitive, user-friendly analytical tools (Gosselin, 2021; Armstrong, 2014; Gómez-López et al., 2019)

In this context, `plotastic` emerges as a tool designed to democratize access to sophisticated statistical analysis, offering a user-centric interface that caters to researchers across varying levels of statistical proficiency. By integrating robust statistical methodologies within an accessible framework, `plotastic` aims to contribute to enhancing the reproducibility and integrity of research in the biosciences (Gomez-Cabrero et al., 2014).

initially, the need to develop `plotastic` arose during this project. The first is to address the author's need for a tool that could handle the complex, multidimensional data generated by e.g. qPCR experiments. These experiments typically involve the analysis of multiple genes across several time points and biological replicates, resulting in datasets that are challenging to analyze manually. The author's experience with traditional statistical software, such as Prism, revealed that these tools required extensive manual input, making them unsuitable for the efficient analysis of complex, multidimensional data. - The second was to increase speed. This is required for developing methods

Since `plotastic` optimizes the analysis of *semi-big data*, we introduce the term *semi-automation* to distinguish itself from the fully automated pipelines used for *big data*. Semi-automation is defined as the following aspects:

1. **Semi-big input:** The input size is oriented towards *semi-big data*, which is characterized as being manageable by manual analysis, yet highly time inefficient, and probably impossible to re-analyse by someone else than the researcher.
2. **Standardized input** The input follows a standardized format (e.g. long-format)
3. **Minimize user configuration:** User configuration is strictly minimized. The user is never asked to pass the same parameters twice. This reduces the risk of human error and time spent on configuration.
4. **Default configuration provides acceptable results:** If the user does not provide any manual configuration, the pipeline should provide acceptable results. Options should be provided to allow a level of flexibility to adapt the pipeline to the user's needs.
5. **Small Reviewable Processing Steps:** The analysis steps are structured into small processes that can be combined to form a complete analysis pipeline. That way, each step can act as a stage for quality control to improve error detection and troubleshooting. For

a statistical analysis, that means the processing steps are separated into 3 steps, those being assumption testing, factor analysis and post-hoc testing.

6. **Isolated Steps:** Processing steps should work independently from another, in the best case only depending on the raw data input. If a processing step depends on the output from other steps, the software should tell the user what exact steps it expects.
7. **Human readable outputs:** Every processing step may provide an output that is not necessarily standardized, but is required to be human readable to ensure reviewability.

Challenges: - Reproducibility crisis? - Data is exploding - Demands for rigorous statistical analysis are increasing - Biologists are not trained in statistics

The demands are rising: (Moreno-Indias et al., 2021)

As laid out in the introduction, one can doubt if a PhD student without coding skills is at its max efficiency.

Why does Biomedicine need plotastic?: - Thorough analysis has become a standard, with assumption testing, omnibus tests and post-hoc analyses for every experiment. - But data is increasing - Example of my data? - The number of dedicated statisticians is limited - The know-how of statistics in biology is limited, for example, Some authors ignored the problem of multiple testing while others used the method uncritically with no rationale or discussion (Perneger, 1998; Armstrong, 2014)

Why did I need plotastic?

Why do biologists need plotastic? - Assays output more data in shorter time, e.g. multiplex qPCR - example: 20 genes, 3 timepoints, 11 biological replicates, (all 3 technical replicates already averaged) -  $20 * 3 * 11 = 660$  data points

this is multidimensional data: 660 data points spread across two dimensions: time and gene

- in manual analysis e.g. in Excel, the user has to manually select the data, copy it, paste it into a new sheet, and then perform the statistical test. In Prism, the user has to select the data, click on the statistical test, and then select the data again. This is not only time-consuming, but also prone to

- Re-Analysis: The user has to repeat the process for every gene and timepoint. This is not only time-consuming, but also prone to errors.

shortly Describe Main Packages in more detail: - seaborn: It multidimensional data - pingouin: It's a statistical package

## Statement of Need

Python's data science ecosystem provides powerful tools for both visualization and statistical testing. However, the transition from exploratory data analysis to hypothesis testing can be cumbersome, requiring users to switch between libraries and adapt to different syntaxes. `seaborn` has become a popular choice for plotting in Python, offering an intuitive interface. Its statistical functionality focuses on descriptive plots and bootstrapped confidence intervals (Waskom, 2021). The library `pingouin` offers an extensive set of statistical tests, but it lacks integration with common plotting capabilities (Vallat, 2018). `statannotations` integrates statistical testing with plot annotations, but uses a complex interface and is limited to pairwise comparisons (Charlier et al., 2022).

`plotastic` addresses this gap by offering a unified environment for plotting and statistical analysis. With an emphasis on user-friendly syntax and integration of familiar `seaborn` parameters, it simplifies the process for users already comfortable with `seaborn`. The library ensures a smooth workflow, from data import to hypothesis testing and visualization.

## Example

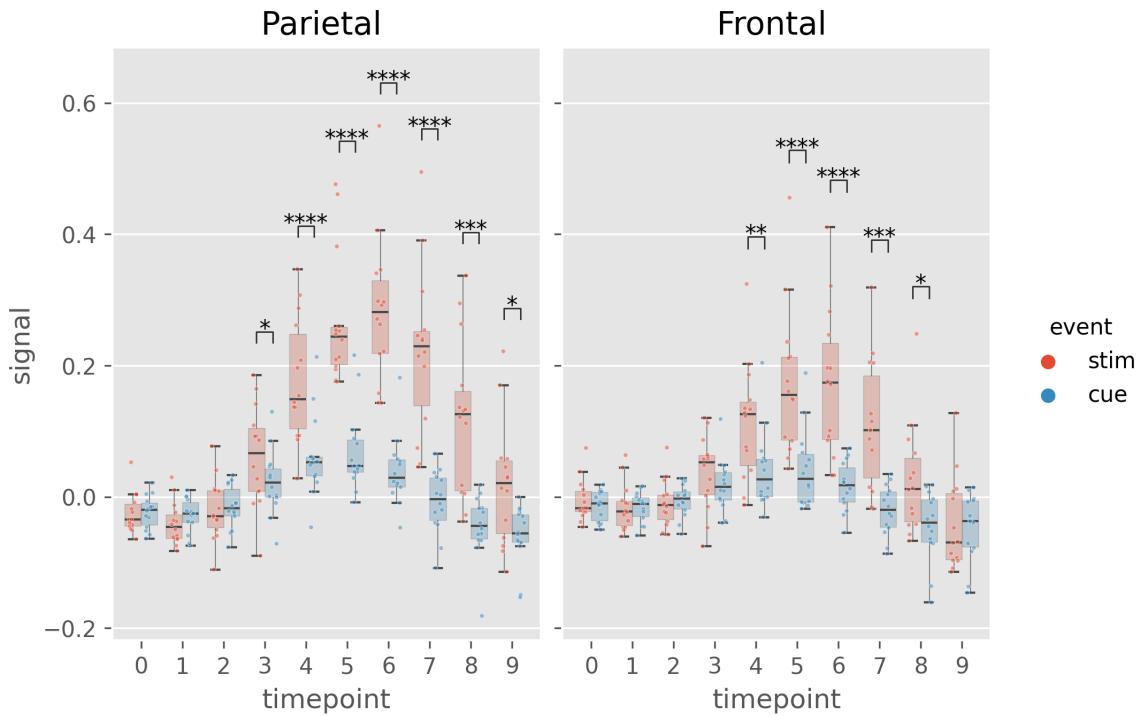
The following code demonstrates how `plotastic` analyzes the example dataset “fmri”, similar to Waskom (2021) (Figure 1).

```
1 ### IMPORT PLOTASTIC
2 import plotastic as plst
3
4 # IMPORT EXAMPLE DATA
5 DF, _dims = plst.load_dataset("fmri", verbose = False)
6
7 # EXPLICITLY DEFINE DIMENSIONS TO FACET BY
8 dims = dict(
9     y = "signal", # y-axis, dependent variable
10    x = "timepoint", # x-axis, independent variable (within-subject factor)
11    hue = "event", # color, independent variable (within-subject factor)
12    col = "region" # axes, grouping variable
13 )
14 # INITIALIZE DATAANALYSIS OBJECT
15 DA = plst.DataAnalysis(
16     data=DF, # Dataframe, long format
17     dims=dims, # Dictionary with y, x, hue, col, row
18     subject="subject", # Datapoints are paired by subject (optional)
19     verbose=False, # Print out info about the Data (optional)
20 )
21 # STATISTICAL TESTS
22 DA.check_normality() # Check Normality
23 DA.check_sphericity() # Check Sphericity
```

```

24 DA.omnibus_rm_anova() # Perform RM-ANOVA
25 DA.test_pairwise()    # Perform Posthoc Analysis
26
27 # PLOTTING
28 (DA
29 .plot_box_strip()    # Pre-built plotting function initializes plot
30 .annotate_pairwise() # Annotate results from DA.test_pairwise()
31     include="__HUE" # Use only significant pairs across each hue
32 )
33

```



**Figure 1:** Example figure of plotastic (version 0.1). Image style was set by plt.style.use('ggplot')

**Table 1:** Results from DA.check\_sphericity(). plotastic assesses sphericity after grouping the data by all grouping dimensions (hue, row, col). For example, DA.check\_sphericity() grouped the ‘fmri’ dataset by “region” (col) and “event” (hue), performing four subsequent sphericity tests for four datasets.

‘region’, ‘event’	spher	W	chi2	dof	pval	group count	n per group
‘frontal’, ‘cue’	True	3.26e+20	-462.7	44	1	10	[14]
‘frontal’, ‘stim’	True	2.45e+17	-392.2	44	1	10	[14]
‘parietal’, ‘cue’	True	1.20e+20	-452.9	44	1	10	[14]
‘parietal’, ‘stim’	True	2.44e+13	-301.9	44	1	10	[14]

**Table 2:** Results of `DA.omnibus_rm_anova()`. `plotastic` performs one two-factor RM-ANOVA per axis (grouping the data by row and col dimensions) using `x` and `hue` as the within-factors. For this example, `DA.omnibus_rm_anova()` grouped the ‘fmri’ dataset by “region” (col), performing two subsequent two-factor RM-ANOVAs. Within-factors are “timepoint” (`x`) and “event” (`hue`). For conciseness, GG-Correction and effect sizes are not shown.

‘region’	Source	SS	ddof1	ddof2	MS	F	p-unc	stars
‘parietal’	timepoint	1.583	9	117	0.175	26.20	3.40e-24	****
‘parietal’	event	0.770	1	13	0.770	85.31	4.48e-07	****
‘parietal’	timepoint * event	0.623	9	117	0.069	29.54	3.26e-26	****
‘frontal’	timepoint	0.686	9	117	0.076	15.98	8.28e-17	****
‘frontal’	event	0.240	1	13	0.240	23.44	3.21e-4	***
‘frontal’	timepoint * event	0.242	9	117	0.026	13.031	3.23e-14	****

## Overview

The functionality of `plotastic` revolves around a seamless integration of statistical analysis and plotting, leveraging the capabilities of `pingouin`, `seaborn`, `matplotlib` and `statannotations` (Vallat, 2018; Waskom, 2021; Hunter, 2007; Charlier et al., 2022). It utilizes long-format `pandas` `DataFrames` as its primary input, aligning with the conventions of `seaborn` and ensuring compatibility with existing data structures (Wickham, 2014; Team, 2020; McKinney, 2010).

`plotastic` was inspired by `seaborn` using the same set of intuitive and consistent parameters (`y`, `x`, `hue`, `row`, `col`) found in each of its plotting functions (Waskom, 2021). These parameters intuitively delineate the data dimensions plotted, yielding ‘faceted’ subplots, each presenting `y` against `x`. This allows for rapid and insightful exploration of multidimensional relationships. `plotastic` extends this principle to statistical analysis by storing these `seaborn` parameters (referred to as dimensions) in a `DataAnalysis` object and intelligently passing them to statistical functions of the `pingouin` library. This approach is based on the impression that most decisions during statistical analysis can be derived from how the user decides to arrange the data in a plot. This approach also prevents code repetition and streamlines statistical analysis. For example, the `subject` keyword is specified only once during `DataAnalysis` initialisation, and `plotastic` selects the appropriate paired or unpaired version of the test. Using `pingouin` alone requires the user to manually pick the correct test and to repeatedly specify the `subject` keyword in each testing function.

In essence, `plotastic` translates plotting parameters into their statistical counterparts. This translation minimizes user input and also ensures a coherent and logical connection between plotting and statistical analysis. The goal is to allow the user to focus on choosing the correct statistical test (e.g. parametric vs. non-parametric) and worry less about specific implementations.

At its core, `plotastic` employs iterators to systematically group data based on various

dimensions, aligning the analysis with the distinct requirements of tests and plots. Normality testing is performed on each individual sample, which is achieved by splitting the data by all grouping dimensions and also the x-axis (hue, row, col, x). Sphericity and homoscedasticity testing is performed on a complete sampleset listed on the x-axis, which is achieved by splitting the data by all grouping dimensions (hue, row, col) (Table 1). For omnibus and posthoc analyses, data is grouped by the row and col dimensions in parallel to the `matplotlib` axes, before performing one two-factor analysis per axis using x and hue as the within/between-factors. (Table 2).

`DataAnalysis` visualizes data through predefined plotting functions designed for drawing multi-layered plots. A notable emphasis within `plotastic` is placed on showcasing individual datapoints alongside aggregated means or medians. In detail, each plotting function initializes a `matplotlib` figure and axes using `plt.subplots()` while returning a `DataAnalysis` object for method chaining. Axes are populated by `seaborn` plotting functions (e.g., `sns.boxplot()`), leveraging automated aggregation and error bar displays. Keyword arguments are passed to these `seaborn` functions, ensuring the same degree of customization. Users can further customize plots by chaining `DataAnalysis` methods or by applying common `matplotlib` code to override `plotastic` settings. Figures are exported using `plt.savefig()`.

`plotastic` also focuses on annotating statistical information within plots, seamlessly incorporating p-values from pairwise comparisons using `statannotations` (Charlier et al., 2022). This integration simplifies the interface and enables options for pair selection in multidimensional plots, enhancing both user experience and interpretability.

For statistics, `plotastic` integrates with the `pingouin` library to support classical assumption and hypothesis testing, covering parametric/non-parametric and paired/non-paired variants. Assumptions such as normality, homoscedasticity, and sphericity are tested. Omnibus tests include two-factor RM-ANOVA, ANOVA, Friedman, and Kruskal-Wallis. Posthoc tests are implemented through `pingouin.pairwise_tests()`, offering (paired) t-tests, Wilcoxon, and Mann-Whitney-U.

To sum up, `plotastic` stands as a unified and user-friendly solution catering to the needs of researchers and data scientists, seamlessly integrating statistical analysis with the power of plotting in Python. It streamlines the workflow, translates `seaborn` parameters into statistical terms, and supports extensive customization options for both analysis and visualization.

## Discussion

Is `plotastic` tested? Coverage? Does it cover every feature? What is not covered

The Architecture of `plotastic` is shown in Appendix A Figure 1.

The full code of an example analysis is shown in section .

Is plotastic USABLE for biologists? - Yes but use is limited by minimal knowledge of Python - However, that is subject to change as Python is becoming more popular in biology and AI assisted coding decreased the barrier to entry significantly. Tools like github copilot are able to generate code, fix bugs and suggest improvements. This is a game changer for biologists that are not familiar with programming. - Furthermore, installing and using plotastic for biologists is overestimated. These steps re needed: - Install anaconda from the internet - Open the terminal - Type `pip install plotastic` - Check Rea

The evaluation of plotastic within this thesis reflects its potential to address key challenges in the field of data analysis. The software integrates a comprehensive suite of statistical tests, such as ANOVA and t-tests, designed for adaptability and ease of use, leveraging the functionalities of pingouin.

In the context of the reproducibility crisis in scientific research, plotastic offers noteworthy contributions, though it is not positioned as a universal remedy. The tool's unique approach to integrating statistical analysis with visual representation establishes a new paradigm, promoting methodological transparency. By mandating that statistical analyses accompany relevant graphical outputs, plotastic ensures that analyses are not only conducted with proper scientific rigor but also documented in a manner that facilitates replication, provided the user possesses proficiency in Python.

Usability is a critical attribute of analytical software, particularly as researchers confront increasingly complex datasets. While the developer's intimate familiarity with plotastic may bias perceptions of its ease of use, it is recognized that novices may initially encounter challenges. Nevertheless, plotastic is distinguished by its user-friendly interface, enabling users with minimal statistical training to perform sophisticated analyses by intuitively mapping plotting concepts to statistical operations.

The transition to a new analytical framework, especially one that incorporates coding, presents a learning curve. However, the advantages of plotastic in terms of analytical clarity, speed, and depth are anticipated to outweigh these initial challenges. Support mechanisms, such as assistance from advanced AI like ChatGPT, are available to mitigate these hurdles, supporting users across varying levels of expertise.

In conclusion, plotastic is posited as a valuable tool in the landscape of scientific research, offering a means to enhance the reproducibility and efficiency of data analysis. Its development ethos emphasizes simplifying complex analytical tasks, thereby contributing to the broader goal of fostering transparent and reproducible research practices.

DO we apply the principles of Semi-Automation to the software?

what features are missing? - Bivariate analysis - Filer to help save the output? - StatResults: System to suggest the correct test, based on the data

# Summarising Discussion

## Time Lapse

Lore Ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lore ipsum dolor sit amet. Lore ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lore ipsum

## Myeloma

Lore Ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lore ipsum dolor sit amet. Lore ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lore ipsum

## Semi-Automated Analysis Improves Agility During Establishing new *in vitro* Methods

Was plotastic useful for me? - Yes incredibly. I was able to perform the statistical tests and visualize the data in a fraction of the time that I would have needed manually. This allowed me to focus on the interpretation of the results and the writing of the manuscript. There was one particular example where my analysis was so fast, that I fed raw datatables during microscopy into python scripts and was able to adapt the experimental technique during the experiment. This allows for an agile and adaptive work environment that is not possible with manual analysis and proved invaluable during development of *in vitro* methods. - These experiments benefited from the use of plotastic, as the

Further research is needed to assess the true impact of semi-automated analysis on the agility

of establishing new *in vitro* methods.

# References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., ... Zheng, X. (2016, March). *TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems* (No. arXiv:1603.04467). arXiv. Retrieved 2024-03-07, from <http://arxiv.org/abs/1603.04467> doi: 10.48550/arXiv.1603.04467
- Alcorta-Sevillano, N., Macías, I., Infante, A., & Rodríguez, C. I. (2020, December). Deciphering the Relevance of Bone ECM Signaling. *Cells*, 9(12), 2630. Retrieved 2023-12-20, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7762413/> doi: 10.3390/cells9122630
- Armstrong, R. A. (2014, September). When to use the Bonferroni correction. *Ophthalmic & Physiological Optics: The Journal of the British College of Ophthalmic Opticians (Optometrists)*, 34(5), 502–508. doi: 10.1111/opo.12131
- Barzilay, R., Ben-Zur, T., Bulvik, S., Melamed, E., & Offen, D. (2009, May). Lentiviral delivery of LMX1a enhances dopaminergic phenotype in differentiated human bone marrow mesenchymal stem cells. *Stem cells and development*, 18(4), 591–601. doi: 10.1089/scd.2008.0138
- Begley, C. G., & Ioannidis, J. P. A. (2015, January). Reproducibility in science: Improving the standard for basic and preclinical research. *Circulation Research*, 116(1), 116–126. doi: 10.1161/CIRCRESAHA.114.303819
- Bianco, P. (2014). "Mesenchymal" stem cells. *Annual review of cell and developmental biology*, 30, 677–704. doi: 10.1146/annurev-cellbio-100913-013132
- Bladé, J., Beksac, M., Caers, J., Jurczyszyn, A., von Lilienfeld-Toal, M., Moreau, P., ... Richardson, P. (2022, March). Extramedullary disease in multiple myeloma: A systematic literature review. *Blood Cancer Journal*, 12(3), 1–10. Retrieved 2023-03-24, from <https://www.nature.com/articles/s41408-022-00643-3> doi: 10.1038/s41408-022-00643-3
- Bondi, A. B. (2000, September). Characteristics of scalability and their impact on performance. In *Proceedings of the 2nd international workshop on Software and performance* (pp. 195–203). New York, NY, USA: Association for Computing Machinery. Retrieved 2024-03-07, from <https://dl.acm.org/doi/10.1145/350391.350432> doi: 10.1145/350391.350432
- Boswell, D., & Foucher, T. (2011). *The Art of Readable Code: Simple and Practical Techniques for Writing Better Code*. "O'Reilly Media, Inc."
- Brooke, J. (1996, January). SUS – a quick and dirty usability scale. In (pp. 189–194).
- Bustin, S. A. (2014, December). The reproducibility of biomedical research: Sleepers awake! *Biomolecular Detection and Quantification*, 2, 35–42. Retrieved 2024-03-18, from <https://www.sciencedirect.com/science/article/pii/S2214753515000030> doi: 10.1016/j.bdq.2015.01.002
- Caplan, A. (1991). Mesenchymal stem cells. *Journal of orthopaedic research : official publication of the Orthopaedic Research Society*, 9(5), 641–50. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/1870029> doi: 10.1002/jor.1100090504
- Caplan, A. I. (1994, July). The mesengenic process. *Clinics in plastic surgery*, 21(3), 429–435.
- Chacon, S., & Straub, B. (2024, March). *Git - Book*. Retrieved 2024-03-07, from <https://git-scm.com/book/de/v2>
- Charlier, F., Weber, M., Izak, D., Harkin, E., Magnus, M., Lalli, J., ... Repplinger, S. (2022, October). *Trevis-md/statannotations: V0.5*. Zenodo. Retrieved 2023-11-16, from <https://zenodo.org/record/7213391> doi: 10.5281/ZENODO.7213391
- Cooper, G. M. (2000). The Cell: A Molecular Approach. 2nd Edition. *Sinauer Associates*, Proliferation in Development and Differentiation. Retrieved from <https://www.ncbi.nlm.nih.gov/books/NBK9906/>
- da Silva Meirelles, L., Chagastelles, P. C., & Nardi, N. B. (2006, June). Mesenchymal stem cells reside in virtually all post-natal organs and tissues. *Journal of cell science*, 119(Pt 11), 2204–2213. doi: 10.1242/jcs.02932

- Dominici, M., Le Blanc, K., Mueller, I., Slaper-Cortenbach, I., Marini, F., Krause, D., ... Horwitz, E. (2006). Minimal criteria for defining multipotent mesenchymal stromal cells. The International Society for Cellular Therapy position statement. *Cytotherapy*, 8(4), 315–317. doi: 10.1080/14653240600855905
- Duvall, P., Matyas, S., & Glover, A. (2007). *Continuous integration: Improving software quality and reducing risk*. Pearson Education. Retrieved from <https://books.google.de/books?id=PV9qfEdv9L0C>
- Ekmekci, B., McAnany, C. E., & Mura, C. (2016, July). An Introduction to Programming for Bioscientists: A Python-Based Primer. *PLOS Computational Biology*, 12(6), e1004867. Retrieved 2024-03-10, from <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004867> doi: 10.1371/journal.pcbi.1004867
- Fazeli, P. K., Horowitz, M. C., MacDougald, O. A., Scheller, E. L., Rodeheffer, M. S., Rosen, C. J., & Klibanski, A. (2013, March). Marrow Fat and Bone—New Perspectives. *The Journal of Clinical Endocrinology and Metabolism*, 98(3), 935–945. Retrieved 2023-12-20, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3590487/> doi: 10.1210/jc.2012-3634
- Friedenstein, A., & Kuralesova, A. I. (1971, August). Osteogenic precursor cells of bone marrow in radiation chimeras. *Transplantation*, 12(2), 99–108.
- Friedenstein, A. J., Piatetzky-Shapiro, I. I., & Petrakova, K. V. (1966, December). Osteogenesis in transplants of bone marrow cells. *Journal of embryology and experimental morphology*, 16(3), 381–390.
- Gabr, M. M., Zakaria, M. M., Refaie, A. F., Ismail, A. M., Abou-El-Mahasen, M. A., Ashamallah, S. A., ... Ghoneim, M. A. (2013). Insulin-producing cells from adult human bone marrow mesenchymal stem cells control streptozotocin-induced diabetes in nude mice. *Cell transplantation*, 22(1), 133–145. doi: 10.3727/096368912X647162
- García-Ortiz, A., Rodríguez-García, Y., Encinas, J., Maroto-Martín, E., Castellano, E., Teixidó, J., & Martínez-López, J. (2021, January). The Role of Tumor Microenvironment in Multiple Myeloma Development and Progression. *Cancers*, 13(2). Retrieved 2021-02-02, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7827690/> doi: 10.3390/cancers13020217
- Glavey, S. V., Naba, A., Manier, S., Clouser, K., Tahri, S., Park, J., ... Ghobrial, I. M. (2017, November). Proteomic characterization of human multiple myeloma bone marrow extracellular matrix. *Leukemia*, 31(11), 2426–2434. Retrieved 2023-09-05, from <https://www.nature.com/articles/leu2017102> doi: 10.1038/leu.2017.102
- Gomez-Cabrero, D., Abugessaisa, I., Maier, D., Teschendorff, A., Merkenschlager, M., Gisel, A., ... Tegnér, J. (2014, March). Data integration in the era of omics: Current and future challenges. *BMC Systems Biology*, 8(2), I1. Retrieved 2024-03-18, from <https://doi.org/10.1186/1752-0509-8-S2-I1> doi: 10.1186/1752-0509-8-S2-I1
- Gómez-López, G., Dopazo, J., Cigudosa, J. C., Valencia, A., & Al-Shahrour, F. (2019, May). Precision medicine needs pioneering clinical bioinformaticians. *Briefings in Bioinformatics*, 20(3), 752–766. doi: 10.1093/bib/bbx144
- Goodman, S. N., Fanelli, D., & Ioannidis, J. P. A. (2016, June). What does research reproducibility mean? *Science Translational Medicine*, 8(341), 341ps12-341ps12. Retrieved 2024-03-18, from <https://www.science.org/doi/10.1126/scitranslmed.aaf5027> doi: 10.1126/scitranslmed.aaf5027
- Gosselin, R.-D. (2021, February). Insufficient transparency of statistical reporting in preclinical research: A scoping review. *Scientific Reports*, 11(1), 3335. Retrieved 2024-03-11, from <https://www.nature.com/articles/s41598-021-83006-5> doi: 10.1038/s41598-021-83006-5
- Gronthos, S., Graves, S. E., Ohta, S., & Simmons, P. J. (1994, December). The STRO-1+ fraction of adult human bone marrow contains the osteogenic precursors. *Blood*, 84(12), 4164–4173.
- Hunter, J. D. (2007, May). Matplotlib: A 2D Graphics Environment. *Computing in Science & Engineering*,

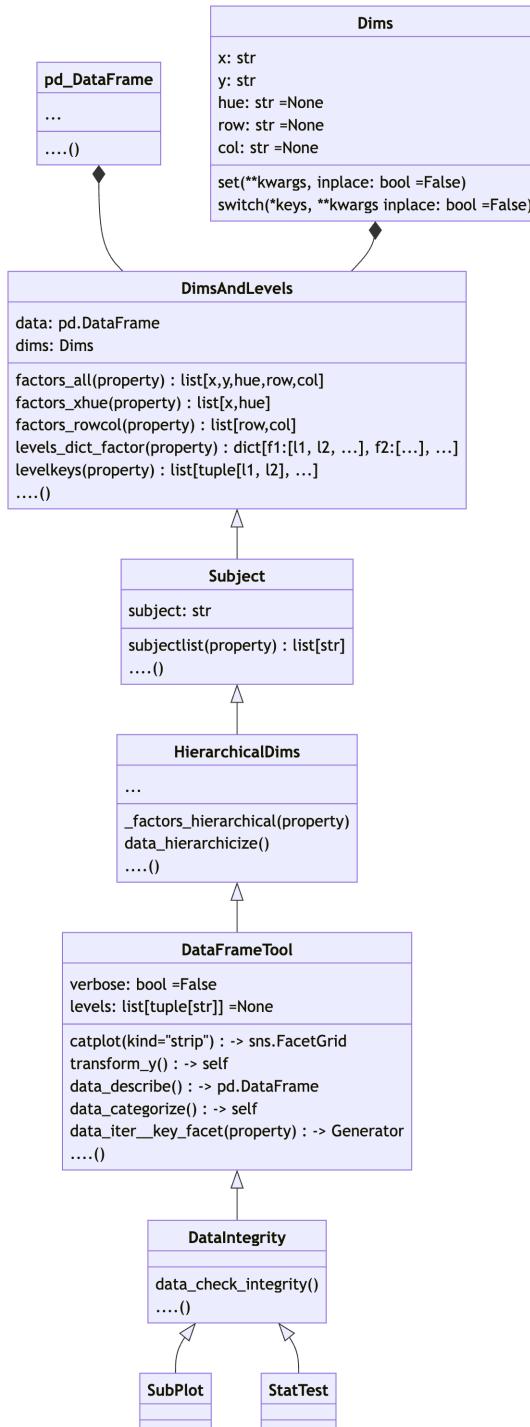
- 9(3), 90–95. Retrieved 2023-11-15, from <https://ieeexplore.ieee.org/document/4160265> doi: 10.1109/MCSE.2007.55
- Incerti, D., Thom, H., Baio, G., & Jansen, J. P. (2019, May). R You Still Using Excel? The Advantages of Modern Software Tools for Health Technology Assessment. *Value in Health*, 22(5), 575–579. Retrieved 2024-03-11, from <https://www.sciencedirect.com/science/article/pii/S1098301519300506> doi: 10.1016/j.jval.2019.01.003
- Jansen, B. J. H., Gilissen, C., Roelofs, H., Schaap-Oziemlak, A., Veltman, J. A., Raymakers, R. A. P., ... Adema, G. J. (2010, April). Functional differences between mesenchymal stem cell populations are reflected by their transcriptome. *Stem cells and development*, 19(4), 481–490. doi: 10.1089/scd.2009.0288
- Kazman, R., Bianco, P., Ivers, J., & Klein, J. (2020, December). *Maintainability* (Report). Carnegie Mellon University. Retrieved 2024-03-07, from <https://kilthub.cmu.edu/articles/report/Maintainability/12954908/1> doi: 10.1184/R1/12954908.v1
- Kibler, C., Schermutzki, F., Waller, H. D., Timpl, R., Müller, C. A., & Klein, G. (1998, June). Adhesive interactions of human multiple myeloma cell lines with different extracellular matrix molecules. *Cell Adhesion and Communication*, 5(4), 307–323. doi: 10.3109/15419069809040300
- Krekel, H., Oliveira, B., Pfannschmidt, R., Bruynooghe, F., Laugher, B., & Bruhin, F. (2004). *Pytest*. Retrieved from <https://github.com/pytest-dev/pytest>
- McKinney, W. (2010, January). Data Structures for Statistical Computing in Python. In (pp. 56–61). doi: 10.25080/Majora-92bf1922-00a
- Moreno-Indias, I., Lahti, L., Nedyalkova, M., Elbere, I., Roshchupkin, G., Adilovic, M., ... Zomer, A. L. (2021, February). Statistical and Machine Learning Techniques in Human Microbiome Studies: Contemporary Challenges and Solutions. *Frontiers in Microbiology*, 12. Retrieved 2024-03-18, from <https://www.frontiersin.org/journals/microbiology/articles/10.3389/fmicb.2021.635781/full> doi: 10.3389/fmicb.2021.635781
- Muruganandan, S., Roman, A. A., & Sinal, C. J. (2009, January). Adipocyte differentiation of bone marrow-derived mesenchymal stem cells: Cross talk with the osteoblastogenic program. *Cellular and molecular life sciences : CMLS*, 66(2), 236–253. doi: 10.1007/s00018-008-8429-z
- Myers, G. J., Sandler, C., & Badgett, T. (2011). *The art of software testing* (3rd ed.). Wiley Publishing. Retrieved from <https://malenezi.github.io/malenezi/SE401/Books/114-the-art-of-software-testing-3-edition.pdf>
- Narzt, W., Pichler, J., Pirklbauer, K., & Zwinz, M. (1998, January). A Reusability Concept for Process Automation Software..
- Nowotschin, S., & Hadjantonakis, A.-K. (2010, August). Cellular dynamics in the early mouse embryo: From axis formation to gastrulation. *Current opinion in genetics & development*, 20(4), 420–427. doi: 10.1016/j.gde.2010.05.008
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., ... Chintala, S. (2019, December). *PyTorch: An Imperative Style, High-Performance Deep Learning Library* (No. arXiv:1912.01703). arXiv. Retrieved 2024-03-07, from <http://arxiv.org/abs/1912.01703> doi: 10.48550/arXiv.1912.01703
- Peng, R. D. (2011, December). Reproducible Research in Computational Science. *Science*, 334(6060), 1226–1227. Retrieved 2024-03-18, from <https://www.science.org/doi/10.1126/science.1213847> doi: 10.1126/science.1213847
- Perneger, T. V. (1998, April). What's wrong with Bonferroni adjustments. *BMJ : British Medical Journal*, 316(7139), 1236–1238. Retrieved 2021-11-24, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1112991/>
- Pittenger, M. F., Mackay, A. M., Beck, S. C., Jaiswal, R. K., Douglas, R., Mosca, J. D., ... Marshak, D. R.

- (1999). Multilineage Potential of Adult Human Mesenchymal Stem Cells. , 284(April), 143–148. doi: 10.1126/science.284.5411.143
- The Python Language Reference.* (n.d.). Retrieved 2024-03-07, from <https://docs.python.org/3/reference/index.html>
- R Core Team. (2018). *R: A language and environment for statistical computing* [Manual]. Vienna, Austria. Retrieved from <https://www.R-project.org/>
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language Models are Unsupervised Multitask Learners.. Retrieved 2024-03-07, from <https://www.semanticscholar.org/paper/Language-Models-are-Unsupervised-Multitask-Learners-Radford-Wu/9405cc0d6169988371b2755e573cc28650d14dfe>
- Rajkumar, S. V., & Kumar, S. (2020, September). Multiple myeloma current treatment algorithms. *Blood Cancer Journal*, 10(9), 94. Retrieved 2023-07-03, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7523011/> doi: 10.1038/s41408-020-00359-2
- Rayhan, A., & Gross, D. (2023). *The Rise of Python: A Survey of Recent Research.* doi: 10.13140/RG.2.2.27388.92809
- Sacchetti, B., Funari, A., Remoli, C., Giannicola, G., Kogler, G., Liedtke, S., ... Bianco, P. (2016). No identical "mesenchymal stem cells" at different times and sites: Human committed progenitors of distinct origin and differentiation potential are incorporated as adventitial cells in microvessels. *Stem Cell Reports*, 6(6), 897–913. Retrieved from <http://dx.doi.org/10.1016/j.stemcr.2016.05.011> doi: 10.1016/j.stemcr.2016.05.011
- Sandve, G. K., Nekrutenko, A., Taylor, J., & Hovig, E. (2013, October). Ten Simple Rules for Reproducible Computational Research. *PLoS Computational Biology*, 9(10), e1003285. Retrieved 2024-03-07, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3812051/> doi: 10.1371/journal.pcbi.1003285
- Shenghui, H., Nakada, D., & Morrison, S. J. (2009). Mechanisms of Stem Cell Self-Renewal. *Annual Review of Cell and Developmental Biology*, 25(1), 377–406. Retrieved from <https://doi.org/10.1146/annurev.cellbio.042308.113248> doi: 10.1146/annurev.cellbio.042308.113248
- Smith, A. M., Niemeyer, K. E., Katz, D. S., Barba, L. A., Githinji, G., Gymrek, M., ... Vanderplas, J. T. (2018). Journal of Open Source Software (JOSS): Design and first-year review. *PeerJ Preprints*, 4, e147. doi: 10.7717/peerj-cs.147
- Stock, P., Bruckner, S., Winkler, S., Dollinger, M. M., & Christ, B. (2014, April). Human bone marrow mesenchymal stem cell-derived hepatocytes improve the mouse liver after acute acetaminophen intoxication by preventing progress of injury. *International journal of molecular sciences*, 15(4), 7004–7028. doi: 10.3390/ijms15047004
- Tam, P. P., & Beddington, R. S. (1987, January). The formation of mesodermal tissues in the mouse embryo during gastrulation and early organogenesis. *Development (Cambridge, England)*, 99(1), 109–126.
- Tanavalee, C., Luksanapruksa, P., & Singhatanadighe, W. (2016, June). Limitations of Using Microsoft Excel Version 2016 (MS Excel 2016) for Statistical Analysis for Medical Research. *Clinical Spine Surgery*, 29(5), 203. Retrieved 2024-03-11, from [https://journals.lww.com/jspinaldisorders/fulltext/2016/06000/limitations\\_of\\_using\\_microsoft\\_excel\\_version\\_2016.5.aspx](https://journals.lww.com/jspinaldisorders/fulltext/2016/06000/limitations_of_using_microsoft_excel_version_2016.5.aspx) doi: 10.1097/BSD.0000000000000382
- Team, T. P. D. (2020, February). *Pandas-dev/pandas: Pandas.* Zenodo. Retrieved from <https://doi.org/10.5281/zenodo.3509134> doi: 10.5281/zenodo.3509134
- Terpos, E., Ntanasis-Stathopoulos, I., Gavriatopoulou, M., & Dimopoulos, M. A. (2018, January). Pathogenesis of bone disease in multiple myeloma: From bench to bedside. *Blood Cancer Journal*, 8(1), 7. doi: 10.1038/s41408-017-0037-4
- Ullah, I., Subbarao, R. B., & Rho, G. J. (2015). Human mesenchymal stem cells - current trends and future prospective Bioscience Reports. doi: 10.1042/BSR20150025

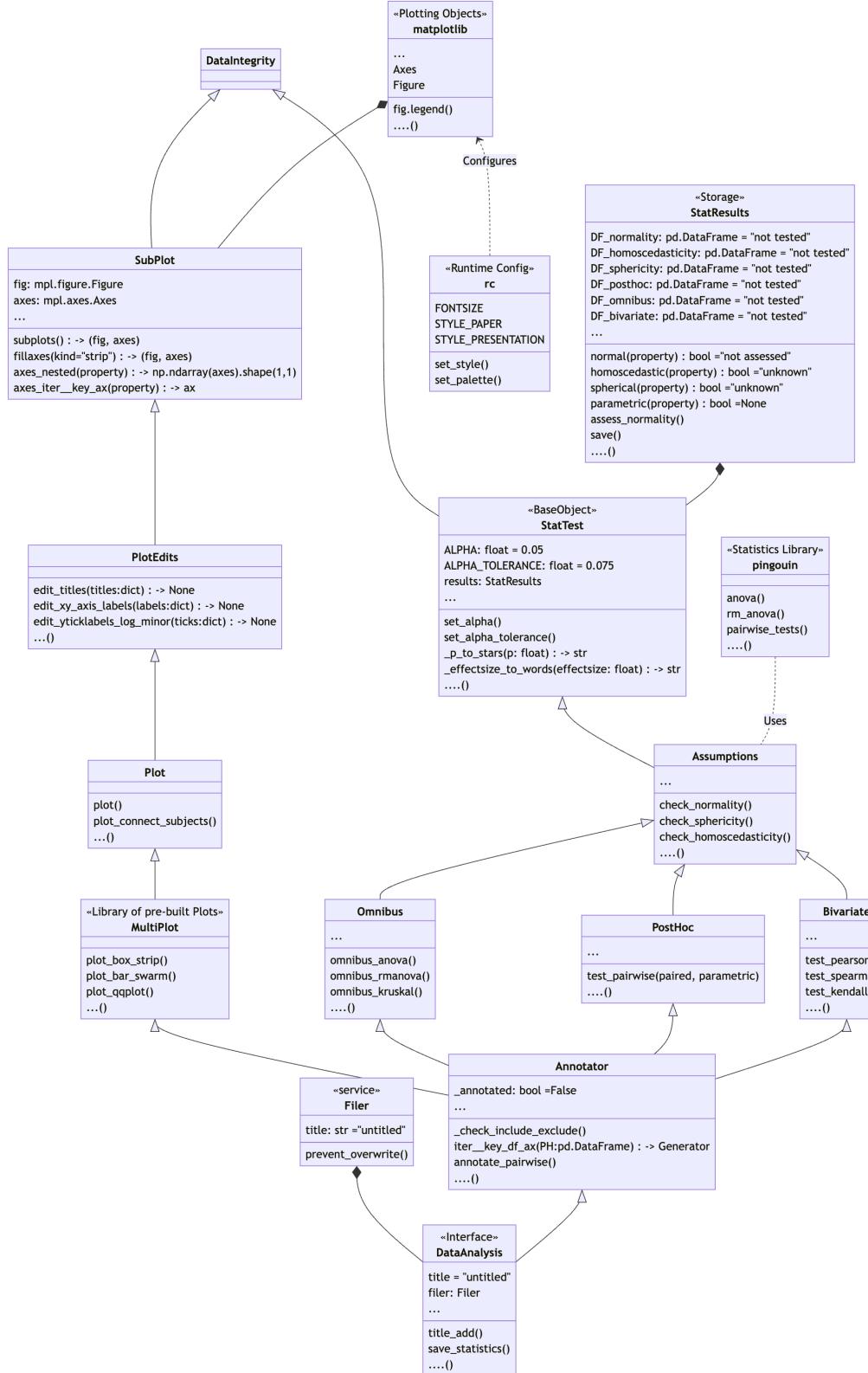
- Vallat, R. (2018, November). Pingouin: Statistics in Python. *Journal of Open Source Software*, 3(31), 1026. Retrieved 2023-05-29, from <https://joss.theoj.org/papers/10.21105/joss.01026> doi: 10.21105/joss.01026
- van Rossum, G., Lehtosalo, J., & Langa, L. (2014). *PEP 484 – Type Hints / peps.python.org*. Retrieved 2024-03-08, from <https://peps.python.org/pep-0484/>
- Viguet-Carrin, S., Garnero, P., & Delmas, P. D. (2006, March). The role of collagen in bone strength. *Osteoporosis International*, 17(3), 319–336. Retrieved 2023-12-20, from <https://doi.org/10.1007/s00198-005-2035-9> doi: 10.1007/s00198-005-2035-9
- Waskom, M. L. (2021, April). Seaborn: Statistical data visualization. *Journal of Open Source Software*, 6(60), 3021. Retrieved 2023-03-26, from <https://joss.theoj.org/papers/10.21105/joss.03021> doi: 10.21105/joss.03021
- Wickham, H. (2014, September). Tidy Data. *Journal of Statistical Software*, 59, 1–23. Retrieved 2023-11-15, from <https://doi.org/10.18637/jss.v059.i10> doi: 10.18637/jss.v059.i10
- Wilkins, A., Kemp, K., Ginty, M., Hares, K., Mallam, E., & Scolding, N. (2009, July). Human bone marrow-derived mesenchymal stem cells secrete brain-derived neurotrophic factor which promotes neuronal survival in vitro. *Stem cell research*, 3(1), 63–70. doi: 10.1016/j.scr.2009.02.006
- Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., ... Mons, B. (2016, March). The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data*, 3(1), 160018. Retrieved 2024-03-18, from <https://www.nature.com/articles/sdata201618> doi: 10.1038/sdata.2016.18
- Xu, W., Zhang, X., Qian, H., Zhu, W., Sun, X., Hu, J., ... Chen, Y. (2004, July). Mesenchymal stem cells from adult human bone marrow differentiate into a cardiomyocyte phenotype in vitro. *Experimental biology and medicine (Maywood, N.J.)*, 229(7), 623–631.
- Yang, A., Troup, M., & Ho, J. W. (2017, July). Scalability and Validation of Big Data Bioinformatics Software. *Computational and Structural Biotechnology Journal*, 15, 379–386. Retrieved 2024-03-07, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5537105/> doi: 10.1016/j.csbj.2017.07.002

## Appendix A

## Class Diagram of plotastic



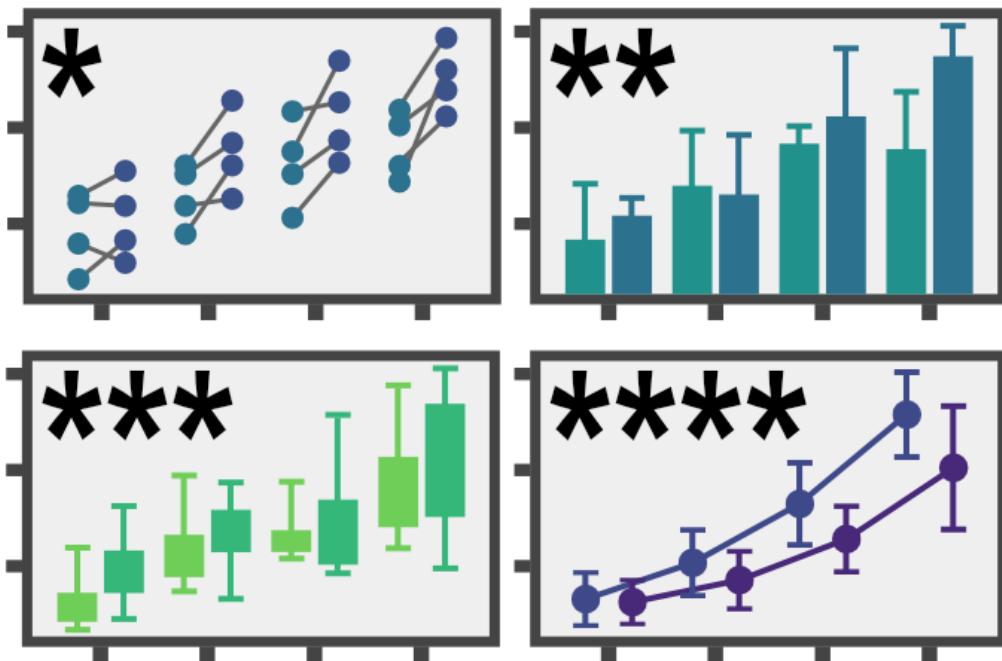
**Appendix Figure 1: Class diagram of plotastic (upper part):** The architecture of plotastic begins with classes that are related to handling a `pandas.DataFrame` object which stores the data, and defining dimensions to group the data (`y`, `x`, `hue`, `col`, `row`). This diagram ends with the classes `SubPlot` and `StatTest` and is continued on the next page. Arrow shapes follow the UML (unified modeling language): A hollow triangle indicates inheritance ("is a") and a filled diamond indicates composition ("has a").



**Appendix Figure 1: (continued)** The architecture of plotastic continues after the class `DataIntegrity` with classes for plotting (`SubPlot`) and statistical testing (`StatTest`) and end with the class `DataAnalysis`, which serves as the main user interface. Arrow shapes follow the UML (unified modeling language): A hollow triangle indicates inheritance ("is a") and a filled diamond indicates composition ("has a").

## Readme of `plotastic`

The following pages are the `README.md` of `plotastic` found in the Python Package Index (PyPi) ([pypi.org/project/plotastic](https://pypi.org/project/plotastic)), and on GitHub ([github.com/markur4/plotastic](https://github.com/markur4/plotastic)).



[code style](#) [black](#)  [codecov](#) [79%](#) [JOSS](#) [10.21105/joss.06304](#)

## plotastic: Bridging Plotting and Statistics

### Installation

#### Install from PyPi:

```
pip install plotastic
```

#### Install from GitHub: (experimental, check CHANGELOG.md)

```
pip install git+https://github.com/markur4/plotastic.git
```

### Requirements

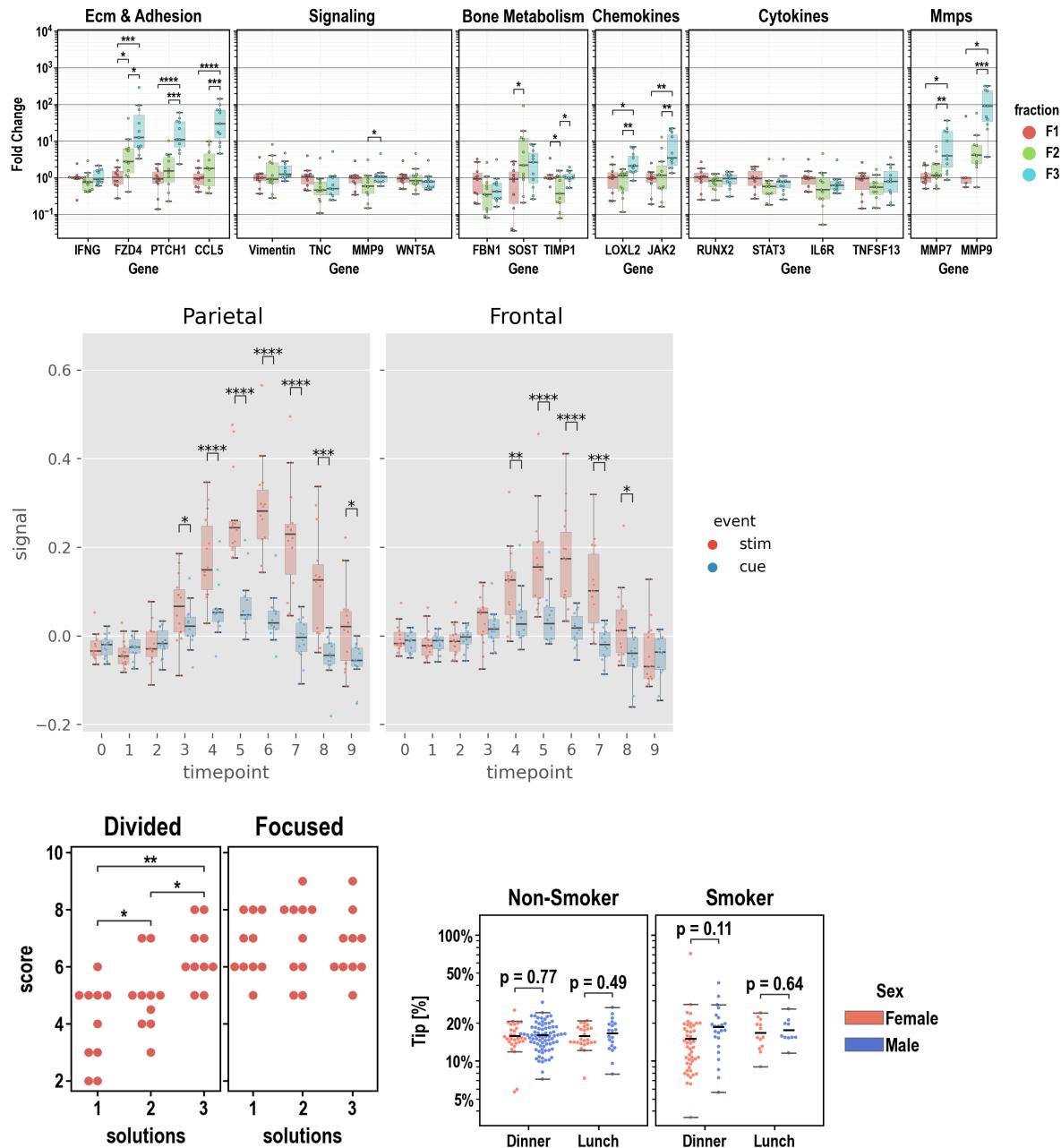
- Python >= 3.11 (*not tested with earlier versions*)

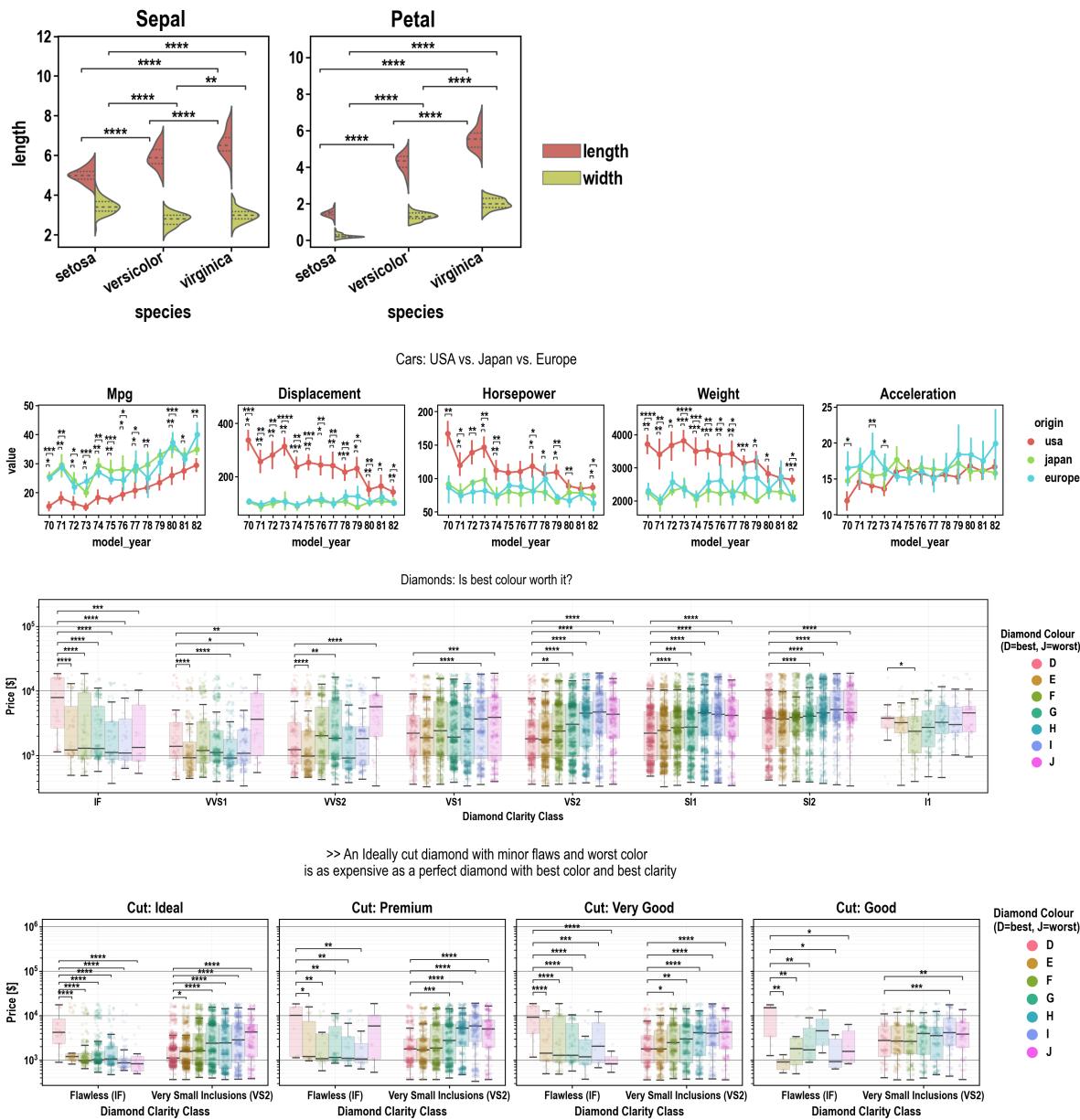
- pandas == 1.5.3 (*pingouin needs this*)
- seaborn <= 0.12.2 (*later versions reworked hue*)

## 📸 Example Gallery

► **(click to unfold)**

🖱️ Click on Images for Code! 🐭





## >About plotastic

### ▶ 🧠 Summary

**plotastic** addresses the challenges of transitioning from exploratory data analysis to hypothesis testing in Python's data science ecosystem. Bridging the gap between **seaborn** and **pingouin**, this library offers a unified environment for plotting and statistical analysis. It simplifies the workflow with a user-friendly syntax and seamless integration with familiar **seaborn** parameters (`y`, `x`, `hue`, `row`, `col`). Inspired by **seaborn**'s consistency, **plotastic** utilizes a **DataAnalysis** object to intelligently pass parameters to **pingouin** statistical functions. The library systematically groups the data according to the needs of statistical tests and plots, conducts visualisation, analyses and supports extensive customization options. In essence, **plotastic** establishes a protocol for configuring statical analyses through plotting parameters.

This approach streamlines the process, translating `seaborn` parameters into statistical terms, providing researchers and data scientists with a cohesive and user-friendly solution in python.!

Workflow:

### 1. Import & Prepare your pandas DataFrame

- We require a long-format pandas dataframe with categorical columns
- If it works with seaborn, it works with plotastic!

### 2. Make a DataAnalysis Object

- `DataAnalysis(DataFrame, dims={x, y, hue, row, col})`
- Check for empty data groups, differing samplesizes, NaN-count, etc. automatically

### 3. Explore Data

- Check Data integrity, unequal samplesizes, empty groups, etc.
- Quick preliminary plotting with e.g. `DataAnalysis.catplot()`

### 4. Adapt Data

- Categorize multiple columns at once
- Transform dependent variable
- Each step warns you, if you introduced NaNs without knowledge!
- etc.

### 5. Perform Statistical Tests

- Check Normality, Homoscedasticity, Sphericity
- Perform Omnibus tests (ANOVA, RMANOVA, Kruskal-Wallis, Friedman)
- Perform PostHoc tests (Tukey, Dunn, Wilcoxon, etc.) based on `pg.pairwise_tests()`

### 6. Plot figure

- Use pre-defined and optimized multi-layered plots with one line (e.g. strip over box)!
- Annotate statistical results (p-values as \*, \*\*, \*\*\*, etc.) with full control over which data to include or exclude!

### 7. Save all results at once!

- One DataAnalysis object holds:
  - One DataFrame in `self.data`
  - One Figure in `self.fig, self.axes`
  - Multiple statistical results: `self.results`
- Use `DataAnalysis.save_statistics()` to save all results to different sheets collected in one .xlsx filesheet per test

## ► Translating Plots into Statistics!

In Principle:

- Categorical data is separable into `seaborn`'s categorization parameters: `x, y, hue, row, col`. We call those "*dimensions*".
- These dimensions are assigned to statistical terms:
  - `y` is the **dependent variable (DV)**
  - `x` and `hue` are **independent variables (IV)** and are treated as **within/between factors** (categorical variables)
  - `row` and `col` are **grouping variables** (categorical variables)
  - A `subject` may be specified for within/paired study designs (categorical variable)

- For each level of **row** or **col** (or for each combination of **row-** and **col** levels), statistical tests will be performed with regards to the two-factors **x** and **hue**

Example with ANOVA:

- Imagine this example data:
  - Each day you measure the tip of a group of people.
  - For each tip, you note down the **day**, **gender**, **age-group** and whether they **smoke** or not.
  - Hence, this data has 4 categorical dimensions, each with 2 or more *levels*:
    - **day**: 4 levels (*monday, tuesday, wednesday, Thursday*)
    - **gender**: 2 levels (*male, female*)
    - **smoker**: 2 levels (*yes, no*)
    - **age-group**: 2 levels (*young, old*)
- Each category is assigned to a place of a plot, and when calling statistical tests, we assign them to statistical terms (in comments):
  -

```
# dims is short for dimensions
dims = dict(      # STATISTICAL TERM:
    y = "tip",      # y-axis, dependent variable
    x = "day",      # x-axis, independent variable (within-
                     subject factor)
    hue = "gender", # color, independent variable (within-
                     subject factor)
    col = "smoker", # axes, grouping variable
    row = "age-group" # axes, grouping variable
)
```

- We perform statistical testing groupwise:
  - For each level-combinations of **smoker** and **age-group**, a two-way ANOVA will be performed (with **day** and **gender** as **between** factors for each datagroup):
    - 1st ANOVA assesses datapoints where **smoker=yes** AND **age-group=young**
    - 2nd ANOVA assesses datapoints where **smoker=yes** AND **age-group=old**
    - 3rd ANOVA assesses datapoints where **smoker=no** AND **age-group=young**
    - 4th ANOVA assesses datapoints where **smoker=no** AND **age-group=old**
  - Three-way ANOVAs are not possible (yet), since that would require setting e.g. **col** as the third factor, or implementing another dimension (e.g. **hue2**).

#### ► ! Disclaimer about Statistics

This software was inspired by ...

- ... ***Intuitive Biostatistics*** - Fourth Edition (2017); Harvey Motulsky
- ... ***Introduction to Statistical Learning with applications in Python*** - First Edition (2023); Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani, Jonathan Taylor
- ... talking to other scientists struggling with statistics

 **plotastic** can help you with...

- ... gaining some practical experience when learning statistics
- ... quickly gain statistical implications about your data without switching to another software
- ... making first steps towards a full statistical analysis
- ... plotting publication grade figures (check statistics results with other software)
- ... publication grade statistical analysis **IF** you really know what you're doing OR you have back-checked your results by a professional statistician
- ... quickly test data transformations (log)

### 🚫 **plotastic** can NOT ...

- ... replace a professional statistician
- ... teach you statistics, you need some basic knowledge (but is awesome for practicing!)
- ... test for multicollinearity (Absence of multicollinearity is required by ANOVA!)
- ... perform stringent correction for multiple testing (e.g. bonferroni), as statistical tests are applied to sub-facets of the whole dataframe for each axes, which depends on the definition of x, hue, col, etc. Hence, corrected p-values might over-estimate the significance of your results.

### 🟡 Be critical and responsible with your statistical analysis!

- **Expect Errors:** Don't trust automated systems like this one!
- **Document your work in ridiculous detail:**
  - Include the applied tests, the number of technical replicates and the number of biological/independent in each figure legend
  - State explicitly what each datapoint represents:
    - 1 datapoint = 1 Technical replicate?
    - 1 datapoint = The mean of all technical replicate per independent replicate/subject?
  - State explicitly what the error-bars mean: Standard deviation? Confidence interval?
  - (Don't mix technical with biological/independent variance)
  - Report if/how you removed outliers
  - Report if you did or did not apply correction methods (multiple comparisons, Greenhouse Geyser, etc.) and what your rationale is (exploratory vs. confirmatory study? Validation through other methods to reduce Type I error?)
- **Check results with professionals:**
  - "Here is my data, here is my question, here is my analysis, here is my interpretation. What do you think?"

### ► ✅ Feature List

- ✅ : Complete and tested
- 👍 : Complete
- 📅 : Planned or unfinished (no date)
- 🧑‍💨 : Maybe..? (Rather not...)
- 🚫 : Not planned, don't want
- 😢 : Help Please..?

### ▼ Plotting

- 👍 Make and Edit Plots: Implemented ✅

- All (non-facetgrid) seaborn plots should work, not tested
- QQ-Plot
- Kaplan-Meyer-Plot
- Interactive Plots (where you click stuff and adjust scale etc.)
  - That's gonna be a lot of work!
- Support for `seaborn.FacetGrid`
  - Why not? - `plotastic` uses `matplotlib` figures and fills its axes with `seaborn` plot functions. In my opinion, that's the best solution that offers the best adaptability of every plot detail while being easy to maintain
- Support for `seaborn.objects` (same as Facetgrid)
  - Why not? - I don't see the need to refactor the code
- NEED HELP WITH: The hidden state of `matplotlib` figures/plots/stuff that gets drawn:
  - I want to save the figure in `DataAnalysis.fig` attribute. As simple as that sounds, `matplotlib` does weird stuff, not applying changes after editing the plot.
  - It'd be cool if I could control the changes to a `DataAnalysis` object better (e.g. using `inplace=True` like with `pd.DataFrame`). But I never figured out how to control `matplotlib` figure generation, even with re-drawing the figure with `canvas`. It's a mess and I wasted so much time already.

#### ▼ Multi-Layered Plotting

- Box-plot + swarm
- Box-plot + strip
- Violin + swarm/strip

#### ▼ Statistics

- Assumption testing
  - Normality (e.g. Shapiro-Wilk)
  - Homoscedasticity (e.g. Levene)
  - Sphericity (e.g. Mauchly)
- Omnibus tests
  - ANOVA, RMANOVA, Kruskal-Wallis, Friedman
  - Mixed ANOVA
  - Annotate Results into Plot
- PostHoc
  - `pg.pairwise_tests()`
    - Works with all primary options. That includes all parametric, non-parametric, paired, unpaired, etc. tests (*t-test*, paired *t-test*, *MWU*, *Wilcoxon*, etc.)
  - Annotate Stars into plots (\*, \*\*, etc.)
    - Specific pairs can be included/excluded from annotation
  - Make correction for multiple testing go over complete DataFrame and not Facet-wise:
- Bivariate
  - Find and Implement system to switch between numerical and categorical x-axis
    - Function to convert numerical data into categorical data by binning?
  - Pearson, Spearman, Kendall

### ▼ Analysis Pipelines

*Idea: Put all those statistical tests into one line. I might work on this only after everything's implemented and working confidently and well!*

- 🐦 `between_samples(parametric=True)`: ANOVA + Tukey (if Normality & Homoscedasticity are given)
- 🐦 `between_samples(parametric=False)`: Kruskal-Wallis + Dunn
- 🐦 `within_samples(parametric=True)`: RM-ANOVA + multiple paired t-tests (if Normality & Sphericity are given)
- 🐦 `within_samples(parametric=False)`: Friedman + multiple Wilcoxon



## How To Use

### Documentations

#### 1. Example Gallery

1. Quick Example: FMRI
2. qPCR (paired, parametric)
3. Cars (unpaired, non-parametric)
4. Diamonds (unpaired, non-parametric)
5. Attention (paired/mixed, parametric)
6. Iris (unpaired, parametric)
7. Tips (unpaired, parametric)

#### 2. Data

1. Set/Switch Dimensions

#### 3. Plotting

1. Quick & Simple: MultiPlots
2. Constructing Plots
3. Legends
4. Styles

### Quick Example

#### Import plotastic and example Data

```
import matplotlib.pyplot as plt
import plotastic as plst

# Import Example Data (Long-Format)
DF, _dims = plst.load_dataset("fmri", verbose = False)
DF.head()
```

Assign each column to a dimension (y, x, hue, col, row):

```

dims = dict(
    y = "signal",      # y-axis, dependent variable
    x = "timepoint",  # x-axis, independent variable & within-subject
    factor
    hue = "event",    # color, grouping variable & within-subject factor
    col = "region"    # axes, grouping variable
)

```

**Initialize DataAnalysis Object**

```

DA = plst.DataAnalysis(
    data=DF,           # Dataframe, long format
    dims=dims,         # Dictionary with y, x, hue, col, row
    subject="subject", # Datapoints are paired by subject (optional)
    verbose=False,     # Print out info about the Data (optional)
)

```

**Perform Statistics**

No arguments need to be passed, although `**kwargs`, are passed to respective `pingouin` functions.

```

DA.check_normality() # Normal Distribution?
DA.check_sphericity() # Sphericity?
DA.omnibus_rm_anova() # Repeated Measures ANOVA
DA.test_pairwise() # Post-hoc tests

```

**Save Results:**

Output is one excel file containing results of all performed tests (normality, anova, t-tests, etc.) in different sheets

```
DA.save_statistics("example.xlsx")
```

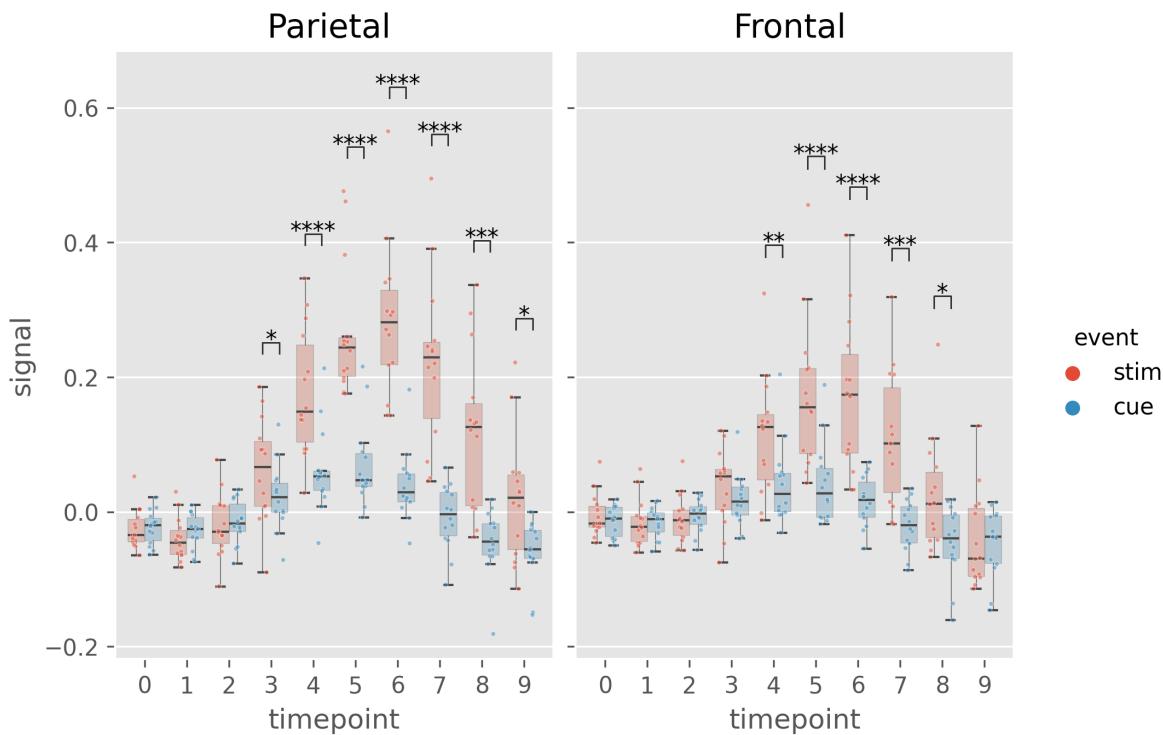
**Annotate post-hoc results into plot:**

```

(DA
    .plot_box_strip() # Pre-built plotting function initializes plot
    .annotate_pairwise( # Annotate results from DA.test_pairwise()
        include="__HUE" # Use only significant pairs across each hue
    )
)

```

```
# Saving the plot like matplotlib!
plt.savefig("example.png", dpi=200, bbox_inches="tight")
```



## 🧪 Testing

▶ (click to unfold)

- Download/Clone repository
- Install development tools `pip install .[dev]`
- Run tests
  - Run `pytest ./tests`
  - To include a coverage report run `pytest ./tests -cov--cov-report=html` and open `./htmlcov/index.html` with your browser.

## 🤝 Community Guidelines

▶ (click to unfold)

When interacting with the community, you must adhere to the [Code of Conduct](#)

### Contribute

I am grateful for [pull requests!](#)

- Make sure to understand the code (e.g. see Class diagram in this Readme)
- Run tests before submitting a pull request

### Reporting Issues & Problems

If you need help, please open an [issue](#) on this repository.

- Please provide a minimal example to reproduce the problem.

## Support

If you need help, please open an [issue](#) on this repository.

## ✍ Cite These!

### ► **(click to unfold)**

Kuric et al., (2024). plotastic: Bridging Plotting and Statistics in Python. *Journal of Open Source Software*, 9(95), 6304, <https://doi.org/10.21105/joss.06304>

Vallat, R. (2018). Pingouin: statistics in Python. *Journal of Open Source Software*, 3(31), 1026, <https://doi.org/10.21105/joss.01026>

Waskom, M. L., (2021). seaborn: statistical data visualization. *Journal of Open Source Software*, 6(60), 3021, <https://doi.org/10.21105/joss.03021>.

```
@article{Kuric2024,
  doi = {10.21105/joss.06304},
  url = {https://doi.org/10.21105/joss.06304},
  year = {2024}, publisher = {The Open Journal},
  volume = {9},
  number = {95},
  pages = {6304},
  author = {Martin Kuric and Regina Ebert},
  title = {plotastic: Bridging Plotting and Statistics in Python},
  journal = {Journal of Open Source Software}
}

@article{Waskom2021,
  doi = {10.21105/joss.03021},
  url = {https://doi.org/10.21105/joss.03021},
  year = {2021},
  publisher = {The Open Journal},
  volume = {6},
  number = {60},
  pages = {3021},
  author = {Michael L. Waskom},
  title = {seaborn: statistical data visualization},
  journal = {Journal of Open Source Software}
}

@article{Vallat2018,
  title = "Pingouin: statistics in Python",
  author = "Vallat, Raphael",
  journal = "The Journal of Open Source Software",
  volume = 3,
```

```
number    = 31,  
pages     = "1026",  
month     = nov,  
year      = 2018  
}
```

## Example Analysis “qpcr” using `plotastic`

The following pages are a jupyter notebook from an example analysis using `plotastic` that's found on GitHub ([github.com/markur4/plotastic](https://github.com/markur4/plotastic)). The Dataset is derived from Chapter 1 qPCR of this thesis, exchanging the original gene names with random ones, while preserving gene classes.

## qPCR

April 9, 2024

```
[ ]: import plotastic as plst
import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd

[ ]: # Set Plot Style
plst.set_style("paper")
# plst.set_palette("hls", verbose=True)
plst.set_palette(["#db5f57", "#91db57", "#57d3db"])

#! You chose this color palette: ['#db5f57', '#91db57', '#57d3db', '#db5f57',
'#91db57', '#57d3db', '#db5f57', '#91db57']

['#db5f57',
 '#91db57',
 '#57d3db',
 '#db5f57',
 '#91db57',
 '#57d3db',
 '#db5f57',
 '#91db57']
```

## 1 Example Analysis: qPCR

**Raw Data:** [https://github.com/markur4/plotastic/tree/main/src/plotastic/example\\_data/data](https://github.com/markur4/plotastic/tree/main/src/plotastic/example_data/data)

**Original Source:** (unpublished)

```
[ ]: # Import Example Data
DF, _dims = plst.load_dataset("qpcr", verbose=False)
dims = dict(
    y="fc",
    x="gene",
    hue="fraction",
    # col= 'method',
    row="class",
)
DA = plst.DataAnalysis(DF, dims, subject="subject", verbose=False)
```

```
[ ]: DA.transform_y("log10", inplace=True) # Log transform
DA.check_normality() # -> Only few groups are not normal -> parametric
```

[ ]:

				W	pval	normal	n
	class	gene	fraction				
Bone Metabolism		F1	FBN1	0.936873	0.518768	True	10
			SOST	0.880395	0.131862	True	10
			TIMP1	0.745494	0.004807	False	9
		F2	FBN1	0.954764	0.705148	True	11
			SOST	0.967810	0.863610	True	11
			TIMP1	0.914325	0.274168	True	11
		F3	FBN1	0.915247	0.281020	True	11
			SOST	0.923112	0.345415	True	11
			TIMP1	0.937230	0.488505	True	11
Chemokines		F1	LOXL2	0.930358	0.451421	True	10
			JAK2	0.897331	0.204749	True	10
		F2	LOXL2	0.874630	0.088876	True	11
			JAK2	0.960025	0.772006	True	11
		F3	LOXL2	0.943678	0.564652	True	11
			JAK2	0.878406	0.099301	True	11
Cytokines		F1	RUNX2	0.947142	0.634825	True	10
			STAT3	0.933422	0.482382	True	10
			IL6R	0.927258	0.421472	True	10
			TNFSF13	0.907481	0.264130	True	10
		F2	RUNX2	0.915611	0.283765	True	11
			STAT3	0.907354	0.226836	True	11
			IL6R	0.985709	0.989621	True	11
			TNFSF13	0.958855	0.757330	True	11
		F3	RUNX2	0.924060	0.353917	True	11
			STAT3	0.932663	0.438418	True	11
			IL6R	0.826181	0.020798	False	11
			TNFSF13	0.970421	0.890746	True	11
ECM & Adhesion		F1	IFNG	0.715267	0.001349	False	10
			FZD4	0.981633	0.973303	True	10
			PTCH1	0.911578	0.292008	True	10
			CCL5	0.969121	0.882582	True	10
		F2	IFNG	0.899109	0.180269	True	11
			FZD4	0.979590	0.963841	True	11
			PTCH1	0.986610	0.990734	True	10
			CCL5	0.925780	0.407685	True	10
		F3	IFNG	0.905665	0.216509	True	11
			FZD4	0.923819	0.351743	True	11
			PTCH1	0.957827	0.744318	True	11
			CCL5	0.940093	0.521596	True	11
MMPs		F1	MMP7	0.955749	0.752957	True	9
			MMP9	0.675286	0.005186	False	5
		F2	MMP7	0.926078	0.372552	True	11

		MMP9	0.971128	0.901100	True	10
	F3	MMP7	0.924886	0.361455	True	11
		MMP9	0.913554	0.268549	True	11
Signaling	F1	Vimentin	0.919696	0.354424	True	10
		TNC	0.928589	0.434161	True	10
		NOTCH1	0.922084	0.374662	True	10
		WNT5A	0.903581	0.239742	True	10
	F2	Vimentin	0.957763	0.743507	True	11
		TNC	0.959813	0.769352	True	11
		NOTCH1	0.977556	0.951045	True	11
		WNT5A	0.937156	0.487661	True	11
	F3	Vimentin	0.910924	0.250109	True	11
		TNC	0.884194	0.117578	True	11
		NOTCH1	0.779982	0.005132	False	11
		WNT5A	0.812114	0.013581	False	11

[ ]: DA.check\_sphericity()

		spher	W	chi2	dof	pval	\
class	fraction						
Bone Metabolism	F1	0	True	0.592922	3.658847	2	0.160506
	F2	0	True	0.703252	3.168356	2	0.205116
	F3	0	True	0.832864	1.645964	2	0.439120
Chemokines	F1	0	True	NaN	NaN	1	1.000000
	F2	0	True	NaN	NaN	1	1.000000
	F3	0	True	NaN	NaN	1	1.000000
Cytokines	F1	0	True	0.629185	3.577934	5	0.614197
	F2	0	False	0.262747	11.657816	5	0.040987
	F3	0	False	0.210032	13.610980	5	0.019012
ECM & Adhesion	F1	0	True	0.486690	5.560987	5	0.354712
	F2	0	True	0.295164	8.202615	5	0.149255
	F3	0	True	0.297080	10.586623	5	0.061736
MMPs	F1	0	True	NaN	NaN	1	1.000000
	F2	0	True	NaN	NaN	1	1.000000
	F3	0	True	NaN	NaN	1	1.000000
Signaling	F1	0	True	0.536227	4.812474	5	0.442437
	F2	0	True	0.554009	5.151113	5	0.400336
	F3	0	False	0.117602	18.669462	5	0.002375
		group	count	n per group			
class	fraction						
Bone Metabolism	F1	0	3	[10, 10, 9]			
	F2	0	3	[11, 11, 11]			
	F3	0	3	[11, 11, 11]			
Chemokines	F1	0	2	[10, 10]			
	F2	0	2	[11, 11]			
	F3	0	2	[11, 11]			

Cytokines	F1	0	4	[10, 10, 10, 10]
	F2	0	4	[11, 11, 11, 11]
	F3	0	4	[11, 11, 11, 11]
ECM & Adhesion	F1	0	4	[10, 10, 10, 10]
	F2	0	4	[10, 11, 11, 10]
	F3	0	4	[11, 11, 11, 11]
MMPs	F1	0	2	[9, 5]
	F2	0	2	[11, 10]
	F3	0	2	[11, 11]
Signaling	F1	0	4	[10, 10, 10, 10]
	F2	0	4	[11, 11, 11, 11]
	F3	0	4	[11, 11, 11, 11]

```
[ ]: # Default is (paired) t-test, and since DA has subject: paired=True
DA.test_pairwise()
```

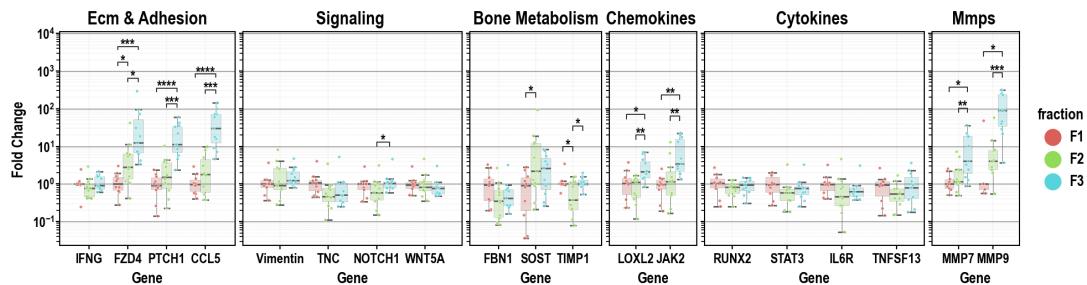
		gene	A	B	mean(A)	\
class	fraction	Contrast				
ECM & Adhesion	-	gene	-	CCL5	FZD4	0.591713
		gene	-	CCL5	IFNG	0.591713
		gene	-	CCL5	PTCH1	0.591713
		gene	-	FZD4	IFNG	0.622994
		gene	-	FZD4	PTCH1	0.622994
..			..	..	..	..
MMPs	NaN	gene * fraction	MMP9	F1	F3	0.256111
		gene * fraction	MMP9	F2	F3	0.677357
	F1	fraction * gene	NaN	MMP7	MMP9	0.032549
	F2	fraction * gene	NaN	MMP7	MMP9	0.185211
	F3	fraction * gene	NaN	MMP7	MMP9	0.742060
			std(A)	mean(B)	std(B)	Paired
class	fraction	Contrast				\
ECM & Adhesion	-	gene	0.253752	0.622994	0.266747	True
		gene	0.253752	-0.026656	0.149430	True
		gene	0.253752	0.469495	0.330886	True
		gene	0.266747	-0.026656	0.149430	True
		gene	0.266747	0.469495	0.330886	True
..			..	..	..	..
MMPs	NaN	gene * fraction	0.802159	1.845550	0.600687	True
		gene * fraction	0.546148	1.845550	0.600687	True
	F1	fraction * gene	0.228544	0.256111	0.802159	True
	F2	fraction * gene	0.361750	0.677357	0.546148	True
	F3	fraction * gene	0.567249	1.845550	0.600687	True
			Parametric	T	dof	\
class	fraction	Contrast				
ECM & Adhesion	-	gene	True	-0.327586	10.0	

		gene	True	7.882620	10.0	
		gene	True	1.320783	10.0	
		gene	True	7.532512	10.0	
		gene	True	2.105924	10.0	
..			..	..	..	
MMPs	NaN	gene * fraction	True	-3.513968	4.0	
		gene * fraction	True	-5.680475	9.0	
	F1	fraction * gene	True	-0.543884	4.0	
	F2	fraction * gene	True	-3.811156	9.0	
	F3	fraction * gene	True	-15.767066	10.0	
			alternative		p-unc	BF10 \
class		fraction Contrast				
ECM & Adhesion -		gene	two-sided	7.499799e-01	0.312	
		gene	two-sided	1.339935e-05	1643.947	
		gene	two-sided	2.160003e-01	0.598	
		gene	two-sided	1.987311e-05	1165.781	
		gene	two-sided	6.146203e-02	1.461	
..			..	..	..	
MMPs	NaN	gene * fraction	two-sided	2.458360e-02	3.686	
		gene * fraction	two-sided	3.016844e-04	111.751	
	F1	fraction * gene	two-sided	6.154168e-01	0.448	
	F2	fraction * gene	two-sided	4.145762e-03	12.636	
	F3	fraction * gene	two-sided	2.163081e-08	4.845e+05	
			hedges	**p-unc	Sign.	\
class		fraction Contrast				
ECM & Adhesion -		gene	-0.115597	ns	False	
		gene	2.856874	****	signif.	
		gene	0.398765	ns	False	
		gene	2.890772	****	signif.	
		gene	0.491362	0.061	toler.	
..			..	..	..	
MMPs	NaN	gene * fraction	-2.179426	*	signif.	
		gene * fraction	-2.024511	***	signif.	
	F1	fraction * gene	-0.255634	ns	False	
	F2	fraction * gene	-0.939861	**	signif.	
	F3	fraction * gene	-1.817138	****	signif.	
				pairs	cross	
class		fraction Contrast				
ECM & Adhesion -		gene		(CCL5, FZD4)	x	
		gene		(CCL5, IFNG)	x	
		gene		(CCL5, PTCH1)	x	
		gene		(FZD4, IFNG)	x	
		gene		(FZD4, PTCH1)	x	
..				..	..	

MMPs	NaN	gene * fraction	((MMP9, F3), (MMP9, F1))	hue
		gene * fraction	((MMP9, F3), (MMP9, F2))	hue
	F1	fraction * gene	((MMP9, F1), (MMP7, F1))	x
	F2	fraction * gene	((MMP9, F2), (MMP7, F2))	x
	F3	fraction * gene	((MMP9, F3), (MMP7, F3))	x

[167 rows x 19 columns]

```
[ ]: # Plot
(
  DA.switch("row", "col", verbose=False)
  .set(y="fc", inplace=False) # set y back to fc to display non-log values
  .plot_box_strip(
    subplot_kws=dict(
      figsize=(10, 2.5),
      width_ratios=[4, 5, 3, 2, 5, 2],
    ),
    strip_kws=dict(alpha=0.8),
  )
  .edit_grid()
  .edit_y_scale_log(10)
  .edit_xy_axis_labels(y_leftmost_col="Fold Change", x="Gene")
  .annotate_pairwise(include="__HUE")
)
plt.savefig("qpcr1.png", dpi=300, bbox_inches="tight")
```



## Appendix B



**Statement of individual author contributions and of legal second publication rights to manuscripts included in the dissertation**

**Manuscript 1: Research Article (submitted, under revision)**

Martin Kuric (MK), Susanne Beck, Doris Schneider, Wyonna Rindt, Marietheres Evers, Jutta Meißner-Weigl, Sabine Zeck, Melanie Krug, Marietta Herrmann, Tanja Nicole Hartmann, Ellen Leich, Maximilian Rudert, Denitsa Docheva, Anja Seckinger, Dirk Hose, Franziska Jundt, Regina Ebert (RE) (2024): Keep it Together: Describing Myeloma Dissemination *in vitro* with hMSC-Interacting Subpopulations and their Aggregation/Detachment Dynamics, **Cancer Research Communications**

Participated in	<b>Author Initials</b> , Responsibility decreasing from left to right				
Study Design	<u>MK</u>	Regina Ebert	Wyonna Rindt		
Methods Development	<u>MK</u>	Doris Schneider			
Data Collection	<u>MK</u>	Doris Schneider			
Data Analysis and Interpretation	<u>MK</u>	Susanne Beck	Regina Ebert		
Manuscript Writing Writing of Introduction Writing of Materials & Methods Writing of Discussion Writing of First Draft	<u>MK</u>	Regina Ebert			

**Explanations:** The content of this publication exceeds the usual scope (~29 pages Supplemental). It includes not only research findings but also survival data and protocols of new, established methods and their validations. The contribution of Martin Kuric was pivotal and predominant in all aspects of this work. Doris Schneider assisted in the experimental procedures. Susanne Beck analyzed the raw data from RNAseq and survival data, which were interpreted, depicted, and summarized by Martin Kuric.

**Manuscript 2: Data Analysis Software (submitted, passed peer-review, under revision)**

Martin Kuric (MK), Regina Ebert (2024): plotastic: Bridging Plotting and Statistics in Python, **Journal of Open Source Software**

Participated in	<b>Author Initials</b> , Responsibility decreasing from left to right				
Idea, Architectural Design	<u>MK</u>				
Software Development Feature Implementation Testing	<u>MK</u>				
Distribution of Software Documentation Version Control (GitHub) Deployment (PyPi)	<u>MK</u>				
Manuscript Writing Writing of Statement of Need Writing of Example Writing of Overview	<u>MK</u>	Regina Ebert			

**Explanations:** The software was entirely created by Martin Kuric, comprising more than 8000 total lines (including ~2000 testable lines) and is comparable in size to a typical web application. The release of this software involved version control using GitHub, packaging and deployment on PyPi. Regina Ebert gave feedback on submitted manuscript.

**Manuscript 3: Research Letter (published)**

Daniela Simone Maichl, Julius Arthur Kirner, Susanne Beck, Wen-Hui Cheng, Melanie Krug, Martin Kuric (MK), Carsten Patrick Ade, Thorsten Bischler, Franz Jakob, Dirk Hose, Anja Seckinger, Regina Ebert & Franziska Jundt (2023): Identification of NOTCH-driven matrisome-associated genes as prognostic indicators of multiple myeloma patient survival, **Blood Cancer Journal 13:134**

<b>Participated in</b>	<b>Author Initials</b> , Responsibility decreasing from left to right				
Study Design Methods Development	Daniela Simone		Franziska Jundt		
Data Collection	Daniela Simone		Franziska Jundt		
Data Analysis and Interpretation	Daniela Simone	Susanne Beck	Franziska Jundt	<u>MK</u>	
Manuscript Writing Writing of Introduction Writing of Materials & Methods Writing of Discussion Writing of First Draft	Daniela Simone		Franziska Jundt	<u>MK</u>	

**Explanations:** This co-authorship is not a chapter in this dissertation. Martin Kuric produced figures of processed but complex-to-visualize data and gave feedback on submitted manuscript.

**Manuscript 4: Research Paper (under peer-review)**

Wyonna Rindt, Melanie Krug, Shuntaro Yamada, Franziska Sennefelder, Louisa Belz, Wen-Hui Cheng, Azeem Muhammad, Martin Kuric (MK), Marietheres Evers, Ellen Leich, Tanja Nicole Hartmann, Ana Rita Pereira, Marietta Herrmann, Jan Hansmann, Mohammed Ahmed Yassin, Kamal Mustafa, Regina Ebert, and Franziska Jundt (2024): A 3D bioreactor model to study osteocyte differentiation and mechanobiology under perfusion and compressive mechanical loading, **Acta Biomaterialia**

<b>Participated in</b>	<b>Author Initials</b> , Responsibility decreasing from left to right				
Study Design Methods Development	Wyonna Rindt	Franziska Jundt		<u>MK</u>	
Data Collection	Wyonna Rindt	Franziska Jundt		<u>MK</u>	
Data Analysis and Interpretation	Wyonna Rindt	Franziska Jundt		<u>MK</u>	
Manuscript Writing Writing of Introduction Writing of Materials & Methods Writing of Discussion Writing of First Draft	Wyonna Rindt	Franziska Jundt		<u>MK</u>	

**Explanations:** This co-authorship is not a chapter in this dissertation. Martin Kuric contributed by counseling during weekly meetings in tight collaboration with Franziska Jundt's group, assisting Wyonna Rindt during laboratory experiments, image analysis and giving feedback on submitted manuscript.

**Manuscript 5: Research Paper (under revision)**

Marietta Herrmann, Jutta Schneidereit, Susanne Wiesner, Martin Kuric (MK), Maximilian Rudert, Martin Lüdemann, Mugdha Srivastava, Norbert Schütze, Regina Ebert, Denitsa Docheva, Franz Jakob (2024): Peripheral blood cells enriched by adhesion to CYR61 are heterogenous myeloid modulators of tissue regeneration with early endothelial progenitor characteristics, **European Cells and Materials**

<b>Participated in</b>	<b>Author Initials</b> , Responsibility decreasing from left to right				
Study Design Methods Development	Marietta Herrmann				
Data Collection	Marietta Herrmann			<u>MK</u>	

Data Analysis and Interpretation	Marietta Herrmann				
Manuscript Writing Writing of Introduction Writing of Materials & Methods Writing of Discussion Writing of First Draft	Marietta Herrmann			<u>MK</u>	

**Explanations:** This co-authorship is not a chapter in this dissertation. Martin Kuric contributed by establishing and measuring large automated microscopy scans of stained cells for quantifying osteogenic differentiation and giving feedback on submitted manuscript.

<b>Manuscript 6: Research Letter (published)</b>					
<b>Participated in</b>	<b>Author Initials, Responsibility decreasing from left to right</b>				
Study Design Methods Development	Marietheres Evers				
Data Collection	Marietheres Evers				
Data Analysis and Interpretation	Marietheres Evers			<u>MK</u>	
Manuscript Writing Writing of Introduction Writing of Materials & Methods Writing of Discussion Writing of First Draft	Marietheres Evers			<u>MK</u>	

**Explanations:** This co-authorship is not a chapter in this dissertation. Martin Kuric contributed by counseling during regular meetings with Ellen Leich's group and giving feedback on submitted manuscript.

If applicable, the doctoral researcher confirms that she/he has obtained permission from both the publishers (copyright) and the co-authors for legal second publication.

The doctoral researcher and the primary supervisor confirm the correctness of the above mentioned assessment.

Würzburg

---

Doctoral Researcher's Name	Date	Place	Signature
----------------------------	------	-------	-----------

Würzburg

---

Primary Supervisor's Name	Date	Place	Signature
---------------------------	------	-------	-----------



**Statement of individual author contributions to figures/tables of manuscripts included in the dissertation**

**Manuscript 1: Research Article (submitted, under revision)**

Martin Kuric (MK), Susanne Beck, Doris Schneider, Wyonna Rindt, Marietheres Evers, Jutta Meißner-Weigl, Sabine Zeck, Melanie Krug, Marietta Herrmann, Tanja Nicole Hartmann, Ellen Leich, Maximilian Rudert, Denitsa Docheva, Anja Seckinger, Dirk Hose, Franziska Jundt, Regina Ebert1 (2024): Keep it Together: Describing Myeloma Dissemination *in vitro* with hMSC-Interacting Subpopulations and their Aggregation/Detachment Dynamics, **Cancer Research Communications**

Figure	<b>Author Initials</b> , Responsibility decreasing from left to right				
1	<u>MK</u>	Doris Schneider			
2	<u>MK</u>	Doris Schneider			
3	<u>MK</u>	Doris Schneider	Sabine Zeck	Wyonna Rindt	Melanie Krug
4	<u>MK</u>	Doris Schneider	Susanne Beck		
5	<u>MK</u>	Susanne Beck			
6	<u>MK</u>	Susanne Beck			
7	<u>MK</u>				
S1	<u>MK</u>	Doris Schneider	Sabine Zeck	Wyonna Rindt	Melanie Krug
S2	<u>MK</u>	Doris Schneider	Marietta Herrmann		
S3	<u>MK</u>	Doris Schneider	Sabine Zeck		
S4	<u>MK</u>				
S5	<u>MK</u>				
S6	<u>MK</u>	Susanne Beck			
Table	<b>Author Initials</b> , Responsibility decreasing from left to right				
1	<u>MK</u>	Susanne Beck			
2	<u>MK</u>	Susanne Beck			
S1	<u>MK</u>	Doris Schneider			
S2	<u>MK</u>	Susanne Beck			
S3	<u>MK</u>				
S4	<u>MK</u>	Doris Schneider			

**Manuscript 2: Data Analysis Software (submitted, passed peer-review, under revision)**

Martin Kuric, Regina Ebert (2024): plotastic: Bridging Plotting and Statistics in Python, **Journal of Open Source Software**

Figure	<b>Author Initials</b> , Responsibility decreasing from left to right				
1	<u>MK</u>				
Table	<b>Author Initials</b> , Responsibility decreasing from left to right				
1	<u>MK</u>				
2	<u>MK</u>				

Documentation	<b>Author Initials</b> , Responsibility decreasing from left to right				
README	<u>MK</u>				
Example Gallery	<u>MK</u>				
Features	<u>MK</u>				
Testing	<b>Author Initials</b> , Responsibility decreasing from left to right				
Test-Code (Pytest)	<u>MK</u>				
Continuous Integration	<u>MK</u>				

**Explanations:** All files are available on GitHub (<https://github.com/markur4/plotastic>) and installable via pypi.com. Documentations are found in the Readme, including example gallery and feature explanation. Software tests was written using pytest. Coverage of code by tests is reviewable with codecov (<https://app.codecov.io/gh/markur4/plotastic>). Continuous Integration is implemented using GitHub actions.

#### Manuscript 3: Research Letter (published)

Daniela Simone Maichl, Julius Arthur Kirner, Susanne Beck, Wen-Hui Cheng, Melanie Krug, Martin Kuric, Carsten Patrick Ade, Thorsten Bischler, Franz Jakob, Dirk Hose, Anja Seckinger, Regina Ebert & Franziska Jundt (2023): Identification of NOTCH-driven matrisome-associated genes as prognostic indicators of multiple myeloma patient survival, **Blood Cancer Journal** **13:134**

Figure	<b>Author Initials</b> , Responsibility decreasing from left to right				
1 a	Daniela Simone			Susanne Beck	
1 b	Daniela Simone			Susanne Beck	
1 c	Daniela Simone			Susanne Beck	<u>MK</u>
1 d	Daniela Simone			Susanne Beck	<u>MK</u>
Table	<b>Author Initials</b> , Responsibility decreasing from left to right				
1	Daniela Simone				

**Explanations:** Martin Kuric plotted multidimensional diagrams using python and fine-adjusted them using professional design software (Affinity Publisher, Serif Ltd).

#### Manuscript 4: Research Paper (under peer-review)

Wyonna Rindt, Melanie Krug, Shuntaro Yamada, Franziska Sennefelder, Louisa Belz, Wen-Hui Cheng, Azeem Muhammad, Martin Kuric (MK), Marietheres Evers, Ellen Leich, Tanja Nicole Hartmann, Ana Rita Pereira, Marietta Hermann, Jan Hansmann, Mohammed Ahmed Yassin, Kamal Mustafa, Regina Ebert, and Franziska Jundt (2024): A 3D bioreactor model to study osteocyte differentiation and mechanobiology under perfusion and compressive mechanical loading, **Acta Biomaterialia**

Figure	<b>Author Initials</b> , Responsibility decreasing from left to right				
1	Wyonna Rindt				<u>MK</u>
2	Wyonna Rindt				
3	Wyonna Rindt				
4	Wyonna Rindt				
5	Wyonna Rindt				<u>MK</u>
6	Wyonna Rindt				<u>MK</u>
7	Wyonna Rindt				<u>MK</u>

**Explanations:** Martin Kuric contributed by counseling on experimental procedures and data analysis, such as quantifying normalized fluorescence intensity of immunohistochemistry and qPCR.

**Manuscript 5: Research Paper (under revision)**

Marietta Herrmann, Jutta Schneidereit, Susanne Wiesner, Martin Kuric (MK), Maximilian Rudert, Martin Lüdemann, Mugdha Srivastava, Norbert Schütze, Regina Ebert, Denitsa Docheva, Franz Jakob (2024): Peripheral blood cells enriched by adhesion to CYR61 are heterogenous myeloid modulators of tissue regeneration with early endothelial progenitor characteristics, **European Cells and Materials**

<b>Figure</b>	<b>Author Initials</b> , Responsibility decreasing from left to right				
1	Marietta Herrmann				
2	Marietta Herrmann				
3	Marietta Herrmann				
4	Marietta Herrmann				
5	Marietta Herrmann				
6	Marietta Herrmann				
7	Marietta Herrmann				<u>MK</u>

**Explanations:** Martin Kuric scanned osteogenically differentiated MSCs in Fig. 7 for quantification of alizarin red staining.

**Manuscript 6: Research Letter (published)**

Marietheres Evers, Martin Schreder, Thorsten Stühmer, Franziska Jundt, Regina Ebert, Tanja Nicole Hartmann, Michael Altenbuchinger, Martina Rudelius, Martin Kuric (MK), Wyonna Darleen Rindt, Torsten Steinbrunn, Christian Langer, Sofia Catalina Heredia-Guerrero, Hermann Einsele, Ralf Christian Bargou, Andreas Rosenwald, Ellen Leich (2023): Prognostic value of extracellular matrix gene mutations and expression in multiple myeloma, **Blood Cancer J.** 13(1):43

<b>Figure</b>	<b>Author Initials</b> , Responsibility decreasing from left to right				
1	Marietheres Evers				
2	Marietheres Evers				

**Explanations:** Martin Kuric contributed indirectly through counseling and feedback on submitted manuscript.

I also confirm my primary supervisor's acceptance.

Doctoral Researcher's Name

Date

Place

Signature

## Affidavit

I hereby confirm that my thesis entitled "Development and Semi-Automated Analysis of an in vitro Model for Myeloma Cells Interacting with Mesenchymal Stromal Cells" is the result of my own work. I did not receive any help or support from commercial consultants. All sources and / or materials applied are listed and specified in the thesis.

Furthermore, I confirm that this thesis has not yet been submitted as part of another examination process neither in identical nor in similar form.

Würzburg  
Place, Date

Signature

## Eidesstattliche Erklärung

Hiermit erkläre ich an Eides statt, die Dissertation "Entwicklung und semi-automatisierte Analyse eines in vitro-Modells für Myelomzellen in Interaktion mit mesenchymalen Stromazellen" eigenständig, d.h. insbesondere selbstständig und ohne Hilfe eines kommerziellen Promotionsberaters, angefertigt und keine anderen als die von mir angegebenen Quellen und Hilfsmittel verwendet zu haben.

Ich erkläre außerdem, dass die Dissertation weder in gleicher noch in ähnlicher Form bereits in einem anderen Prüfungsverfahren vorgelegen hat.

Würzburg  
Ort, Datum

Unterschrift

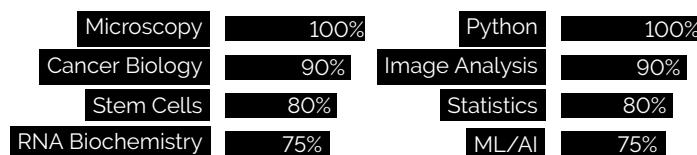
# MARTIN KURIC

Cell Biologist | Data Scientist



## WHO AM I?

As a cancer cell biologist with a strong passion for data analysis and machine learning, I am seeking a position where I can utilize my creativity to automate tasks, solve complex problems or handle big data.



## SELECTED PROJECTS

2024	<b>Python Software “ImageP”</b> Accelerates batch processing of images of different sizes and types by >100%. <code>numpy</code> / <code>skimage</code> / <code>scipy</code>	<a href="#">GitHub Repository</a>
2020-2024	<b>Python Software “plotastic”</b> Published a statistical library that self-configures based on intuitive plotting parameters. <code>pandas</code> / <code>matplotlib</code> / <code>pingouin</code> / <code>seaborn</code>	<a href="#">Journal of Open Source Software</a>   <a href="#">GitHub Repository</a>
2018-2024	<b>Cancer Research Project</b> Worked in a team with up to three technical assistants and published a list of genes with relevance for survival of myeloma patients (under peer-review). <code>Time-Lapse Microscopy</code> / <code>RNAseq</code> / <code>Analysis of Patient Survival</code>	<a href="#">Journal: Cancer Research Communications</a>
26.05.2022	<b>Deep-Learning Assisted Image Cytometry</b> Measurement of per-cell parameters from large automated microscopy scans. <code>Convolutional Neural Networks</code> / <code>Image Segmentation</code>	<a href="#">Poster at "Achilles Conference"</a>

## EDUCATION

28.01.2019 – 2024	<b>Dr. rer. nat. in Biomedicine</b> Research focus: Dissemination of multiple myeloma & mesenchymal stromal cell interactions	<b>Prof. Dr. Regina Ebert</b>   University of Würzburg
01.04.2017 – 2024 parallel to M.Sc. & PhD	<b>Elite Biological Physics</b> Interdisciplinary & international study program for exceptional students of physics or biology.	<b>University of Bayreuth</b>
01.10.15 – 15.08.18	<b>M.Sc. in Biochemistry &amp; Molecular Biology</b> Research focus: RNA biochemistry, small RNAseq, stem cells & piRNAs in <i>S. mediterranea</i>	<b>Prof. Dr. Claus-D. Kuhn</b>   University of Bayreuth
01.10.12 – 14.12.15	<b>B.Sc. in Biochemistry</b> Research focus: Cell biology, mitochondrial inheritance in <i>S. cerevisiae</i>	<b>Prof. Dr. Benedikt Westermann</b>   University of Bayreuth

## LANGUAGES

German, English - C2  
 Slovakian - passive  
 French, Spanish - A2

## SOFT SKILLS

Quality Management  
 Project Management  
 Violent Free Communication

## HOBBIES

Coding - Python  
 Music - Piano & Guitar  
 Gym - Lift. Grow. Repeat

Würzburg

05.03.2024

Location

Date

Martin Kuric

Signature