

# MONTE CARLO METHOD FOR CALCULATING OXYGEN ABUNDANCES AND THEIR UNCERTAINTIES FROM STRONG-LINE FLUX MEASUREMENTS

FEDERICA B. BIANCO<sup>1</sup>, MARYAM MODJAZ<sup>1</sup>, SEUNG MAN OH<sup>1,2</sup>, DAVID FIERROZ<sup>1</sup>, YUQIAN LIU<sup>1</sup>, LISA KEWLEY<sup>3,4</sup>

*Draft version April 24, 2015*

## ABSTRACT

We present the open-source Python code MCZ.py for the determination of the strong-emission-line estimators of oxygen abundance in the standard scales, based on the original IDL-code in Kewley & Dopita (2002), and expanded to include more recently developed scales. The standard strong-line oxygen diagnostics have been used in many areas of astrophysics, including galaxy evolution and supernova (SN) host galaxy studies, to estimate the oxygen abundance in the interstellar medium through various emission line ratios. Oxygen abundance is considered a proxy for total metal abundance, since Oxygen is the most abundant metal. Here we introduce a Python implementation of these methods that through Monte Carlo resampling better characterizes the statistical reddening-corrected oxygen abundance confidence region. Given line flux measurements and their errors, our code produces synthetic Monte Carlo distributions for the oxygen abundance simultaneously in a large number of currently-used metallicity scales in up to 33 scales and subscales, as well as for E(B-V), when possible, and estimates the median values and 16th and 84th percentile confidence regions, for each metallicity diagnostic and for the reddening parameter E(B-V). In addition, it has the option of outputting the full MC parameter distributions, as well as their kernel density estimates. We test our code on emission lines measurements from a sample of supernova host galaxies (with  $z < 0.15$ ) and compare our metallicity results with those from previous methods. We show that our metallicity estimates are consistent with previous methods but yield smaller uncertainties. We also offer visualization tools to assess the spread of the oxygen abundance in the different scales, as well as the shape of the estimated oxygen abundance distribution in each scale, and develop robust metrics for determining the appropriate MC sample size. The code is open access and open source, and it can be found at [www.github.com/nyusngroup/MC\\_Metallicity](http://www.github.com/nyusngroup/MC_Metallicity).

*Subject headings:*

## 1. INTRODUCTION

Small amounts of carbon, oxygen, nitrogen, sulfur and iron and other elements provide a splash of color to the otherwise dominating greyscale of hydrogen and helium in the stars and gas of galaxies. Nevertheless, even this minute presence of heavy elements (all elements heavier than H and He, also called metals or collectively metallicity) is important for many areas of astrophysics. For example, Johnson & Li (2012), amongst others, suggest that if it was not for the relatively high metallicity level in our Solar System, planet formation may not have been possible. With  $Z$  representing the mass fraction of metals, for our own Sun the value is measured to be  $Z=0.0153$  (Caffau et al. 2011), though there suggestions for a lower solar metallicity of  $Z = 0.0134$  (Asplund et al. 2009; Grevesse et al. 2010). Furthermore, when properly observed and estimated, metallicity measurements of galaxies can tightly constrain models of galaxy formation and evolution (e.g., Kewley & Ellison 2008 and references therein), as well as shed light on the metallicity dependence and production conditions for different types of Supernovae (SNe) and long-duration Gamma-

Ray Bursts (GRBs) (e.g., Modjaz et al. 2008; Levesque et al. 2010; Anderson et al. 2010; Modjaz et al. 2011; Kelly & Kirshner 2012; Sanders et al. 2012; Lunnan et al. 2014; Leloudas et al. 2014; Pan et al. 2014).

Most metals are produced in the cores of massive stars during their fusion life cycle but also during the extreme conditions of stellar explosions. For example, the majority of iron is synthesized in thermonuclear explosions (SNe Ia) while nearly all of oxygen and other  $\alpha$ -elements are released in various kinds of core collapse SNe (SNe II and stripped-envelope core-collapse SNe). Since new stars are born from the clouds these explosions enrich, metallicity will increase with each passing generation of stars.

However, for almost all astronomical objects, metallicity cannot be measured directly. The oxygen abundance in the gas-phase is the canonical choice of metallicity indicator for interstellar medium (ISM) studies, since oxygen is the most abundant metal and only weakly depleted onto dust grains (in contrast to refractory elements such as Mg, Si, Fe, with Fe, being depleted by more than a factor of 10 in Orion; see Simón-Díaz & Stasińska 2011). The oxygen abundance<sup>5</sup> is expressed as  $12 + \log_{10}(\frac{O}{H})$ , where  $O$  and  $H$  represent the number of Oxygen and Hydrogen atoms, respectively. In particular, Caffau et al. 2011 measure a solar oxygen abundance of  $12 + \log_{10}(\frac{O}{H}) = 8.76 \pm 0.07$ , while Asplund et al. (2009)

<sup>1</sup> Center for Cosmology and Particle Physics, New York University, 4 Washington Place, New York, NY 10003, USA

<sup>2</sup> NYU Abu Dhabi PO Box 129188 Abu Dhabi, UAE

<sup>3</sup> Australian National University, Research School for Astronomy & Astrophysics, Mount Stromlo Observatory, Cotter Road, Weston, ACT 2611, Australia

<sup>4</sup> Institute of Astronomy, University of Hawaii, 2680 Woodlawn Drive, Honolulu, HI 96822, USA

<sup>5</sup> We note that in many cases in the literature, including here, the terms metallicity and oxygen abundance are used interchangeably.

suggest  $12 + \log_{10}(\frac{O}{H}) = 8.69$ .

Importantly, oxygen exhibits very strong nebular lines in the optical wavelength range in spectra of HII regions (e.g., Pagel et al. 1979; Osterbrock 1989; Tremonti et al. 2004), which can be measured. Thus, many different diagnostic techniques, relying on different lines of oxygen, hydrogen and other elements, have been developed (e.g., Kewley & Dopita 2002; Pettini & Pagel 2004; Kobulnicky & Kewley 2004; Kewley & Ellison 2008), which are discussed in the next section. Ultimately, the purpose of this paper is to present a public code that computes the metallicity from strong emission line fluxes according to the many different abundance diagnostics as well as the associated statistical uncertainties due to the measured emission line flux uncertainties.

### 1.1. The different oxygen abundance diagnostics

Here we present a brief overview of the various observational methods for measuring the gas-phase oxygen abundance - however, for a detailed discussion, and to understand the many caveats, we encourage the reader to read the reviews by e.g. Stasińska (2002); Kewley & Ellison (2008); Moustakas et al. (2010); Stasińska (2010); Dopita et al. (2013); Blanc et al. (2015). The so-called “classical” way to estimate the oxygen abundance is the electron temperature ( $T_e$ ) method, which estimates the electron temperature and density of a nebula using a number of oxygen lines with different ionization states, including the auroral [OIII]  $\lambda 4363$  line, to then directly estimate the OII and OIII abundances and finally obtain the total oxygen abundance, after correcting for the unseen stages of ionization. However, the auroral [OIII]  $\lambda 4363$  line is very weak, except in low-metallicity environments, and saturates at higher metallicity (since at higher metallicities the cooling is dominated by the oxygen NIR fine structure lines) – thus, other methods had to be developed that use other, stronger lines in the spectra of HII regions. These are called *strong-line methods* and are the subject of this manuscript. Strong-line methods can be categorized into two types: theoretical methods, that rely on calibrating various observed line ratios using photoionization methods (basically theoretically simulating HII regions, using stellar model atmospheres, stellar population synthesis and photoionization models) and empirical ones, that calibrate various observed strong line ratios using observed  $T_e$ -based metallicities. While historically there have been large systematic offsets between the  $T_e$  method and the strong line methods, Dopita et al. (2013) demonstrated that the  $T_e$  method gives the same results as the strong line methods, if the energy distribution of the electrons in the HII regions is assumed to not be a simple Maxwell-Boltzmann distribution (as assumed in prior works), but a more realistic  $\kappa$  distribution, as observed in solar system astrophysical plasma. They also find that the effect of changing the specific  $\kappa$  distribution on the strong-line methods is minor.

For theoretical strong-line method, a ratio of oxygen line fluxes to  $H\beta$ , referred to as  $R_{23}$ , is commonly used to determine the metallicity of galaxies (Pagel et al. 1979):

$$R_{23} = \frac{[\text{OII}]\lambda 3727 + [\text{OIII}]\lambda 4959, \lambda 5007}{H\beta},$$

where [OIII] $\lambda 4959, \lambda 5007$  stands for the sum of the two

[OIII] lines. The drawback of this method is that is double-valued with metallicity: the same  $R_{23}$  value may correspond to two metallicity values, and thus other line ratios need to be used to break the degeneracy between the high values (“upper branch”) and the low values (“lower branch”) of the  $R_{23}$  metallicities (e.g., Kewley & Ellison 2008, Moustakas et al. 2010). Furthermore, Kewley & Dopita (2002) showed the importance of ionization parameter, which can be physically understood as corresponding to the maximum velocity of an ionized front that can be driven by the local radiation field of hot massive stars that are ionizing the ISM gas. This ionization parameter needs to be taken into account in the various strong-line methods, as HII regions at the same metallicity but with different ionization parameters produce different line strengths. Calibrations of  $R_{23}$  by McGaugh (1991) (hereafter M91), by Kewley & Dopita (2002) (hereafter KD02), and by Dopita et al. (2013) (hereafter D13) use different theoretical photoionization models and take the ionization parameter into account, while other calibrations such as that of Zaritsky et al. (1994) (hereafter Z94) do not. Thus, Z94 is mostly valid for only metal-rich galaxies. M91 and KD02 use an iterative process to break the  $R_{23}$  degeneracy (KD02 uses different ratios [NII]/[OII] and [NII]/ $H\alpha$ ) and to also constrain the ionization parameter  $q$  in order to arrive at the metallicity estimate.

As to empirical strong-line methods, the most commonly used ones are those by Pettini & Pagel (2004) (hereafter PP04), Pilyugin & Thuan (2005) (hereafter P05) and Marino et al. (2013) (hereafter M13). PP04 used HII regions with measured  $T_e$ -based metallicities to derive empirical fits to strong-line ratios, and introduced the line ratios of  $\frac{[\text{NII}]}{H\beta}$  ( $N2$ ) and  $\frac{[\text{OIII}]}{H\beta} / \frac{[\text{NII}]}{H\alpha}$  ( $O3N2$ ) as metallicity diagnostics. Since PP04.N2 employs two closely spaced lines ( $H\alpha$  and NII), which are not affected by stellar absorption, nor (uncertain) reddening, and are easily observed in one simple spectroscopic setup, it has become a popular scale to use, at least for low- $z$  SN host galaxy studies (see meta-analysis by e.g., Sanders et al. 2012; Modjaz 2012; Leloudas et al. 2014). However, it is important to remember that this scale has a number of short-comings: it does not take into account the impact of the ionization parameter, it was initially derived based on only 137 extragalactic HII regions, and the nitrogen emission line employed saturates at high metallicity, and thus the PP04.N2 method saturates for high-metallicity galaxies (at  $12 + \log_{10}(\frac{O}{H}) > 8.8$ , Kewley & Ellison 2008). An updated calibration by Marino et al. (2013) (whose scale our code also outputs, see below) based on many more  $T_e$ -based metallicities (a sample almost three times larger than that of PP04) derives a significantly shallower slope between  $O3N2$  index and oxygen abundance than the PP04 calibration. In addition, most recently, Berg et al. (2015) suggested that the auroral [OIII]  $\lambda 4363$  line, commonly used for  $T_e$  measurements, is the most problematic auroral line to use amongst those of [OII], [OIII] [NII], [SII], [SIII], giving rise to temperature discrepancies.

As it can be seen, each scale has different advantages and disadvantages and should be used in different metallicity regimes (see detailed discussion in e.g., Kewley & Dopita 2002; Stasińska 2002; Kewley & Ellison 2008;

Moustakas et al. 2010; López-Sánchez et al. 2012; Dopita et al. 2013; Blanc et al. 2015). Thus, this open-source code outputs the oxygen abundance in the main (13 as of version v1.0, Spring 2015) metallicity scales (with the KD02 diagnostic having four outputs and the PP04, P10, M13, M08, and D13 diagnostic having multiple outputs as well). While there is a long-standing debate about which diagnostic to use, as there are systematic metallicity offsets between different methods (recombination lines vs. strong-line method vs. “direct”  $T_e$  method, see the above sources), **the relative metallicity trends can be considered robust, if the analysis is performed self-consistently in the same scale, and trends are seen across different scales (Kewley & Ellison 2008; Moustakas et al. 2010)**. Thus, it is then necessary to obtain statistical error bars for the relative comparisons to be meaningful. Note however, that while there are conversion values between different scales (Kewley & Ellison 2008), they apply for large data sets, since those conversion values were derived based on ten thousands of SDSS galaxies, and thus should be used with caution (or not at all) for smaller samples. In addition, one should note that on account of the ongoing debate about the value of the solar oxygen abundance (Asplund et al. 2009; Caffau et al. 2011), that the absolute oxygen calibration is still uncertain.

Here we introduce the open-source Python package named **MCZ**. This Python code allows the user to quickly produce metallicity values with sensible confidence regions for several metallicity scales at once, given an input set of spectral line measurements and their errors. While we do not mean to advocate for a particular metallicity scale to be adopted, the comparison of multiple scale outputs, and the shape of each metallicity distribution, can guide the user in understanding the reliability of a metallicity estimate, given a set of line fluxes. In § 2 we describe our method, and the input and output values of the code. We furthermore develop robust metrics for determining the appropriate MC sample size and perform benchmark calculations that includes computing the computer memory usage for many different scales. In § 3, we compare our method of obtaining abundance uncertainties to previous methods in the literature and conclude with § 4.

## 2. DESCRIPTION OF METALLICITY CODE

For computing oxygen abundances, we start with the iterative IDL code by Kewley & Dopita (2002), hereafter referred to as IDLKD02, which has been updated in Kewley & Ellison (2008) **LISA: YOUR INPUT HERE: what is the update if any or is it exactly as in Kewley & Ellison 08??**. This code was initially written in IDL. We translated the code into Python and added new, more recent scales (see Section 2.1) and new features, of which the most important is the capability of obtaining uncertainties on the metallicity outputs via Monte Carlo resampling. We release our open-source code on GitHub, as we explain below.

### 2.1. Input and Output of code

The input of the code is a set of spectral emission line fluxes. We assume that the observed emission lines to be used to indicate metallicity originate in HII regions and are not due to non-thermal excitation by, for

example, AGN, interstellar shocks from SNe, or stellar winds. Tests to exclude data contaminated by such non-thermal sources should be executed *prior to running this code* using the recommended line ratios by e.g., Baldwin et al. (1981); Kauffmann et al. (2003); Kewley et al. (2006). Furthermore, these lines should have all the correct flux calibration (at least correct relative calibration) and *should have a signal-to-noise ratio ( $S/N$ ) of at least 3*. The latter is important for the success of the Monte Carlo resampling technique, for the following reason. Synthetic line flux measurements are drawn from within a Gaussian distribution with standard deviation equal to the measurement error, and centered on the measured flux value, as described in detail in Section 2.2. Thus a  $S/N \geq 3$  assures that fewer than  $\sim 1\%$  of the measurement fall below zero (and are thus invalid). The code will check each line  $S/N$  and if the  $S/N \geq 3$  condition is not satisfied for any line, a warning message is issued.

Emission line flux values are fed into our Python implementation as in the original IDLKD02 code. The inputs are emission line flux values, and their uncertainties, for the following lines:  $H\alpha$ ,  $H\beta$ , [OI] 6300, [OII] 3727, [OIII] 4959, [OIII] 5007, [NII] 6584, [SII] 6717, [SII] 6731, [SIII] 9096 and [SIII] 9532 can be used to calculate  $S_{23}$ , but are not often observed since they are in the NIR. Only one metallicity scale based on  $S_{23}$  is implemented in the current version of the code (DP00 from Díaz & Pérez-Montero 2000). The line fluxes are to be stored in an ASCII file, and the measurement errors in a separate ASCII files (consult the README.md<sup>6</sup> in the GitHub repository for details about the input format, and find example files in the repository). If the fluxes for the specified lines are not available, the entry should be set to *NaN*. The oxygen abundance will be calculated only for metallicity scales that use valid, non-*NaN*, line fluxes. If the line fluxes necessary for specific scales are not provided, the output metallicities will default to *NaN*. In absence of measurement errors, the flux errors entry should be set to 0.0 in the input ASCII file. If the errors in the measurements are not provided, the code will specify that it cannot create a measurement distribution and determine a confidence interval, but it will calculate and output the nominal metallicity.

The inputted line fluxes are corrected for reddening by using the observed Balmer decrement, for which  $H\alpha$  and  $H\beta$  flux values need to be provided. We assume case B recombination, and thus the standard value of 2.86 as the intrinsic  $H\alpha/H\beta$  ratio (Osterbrock 1989), and apply the standard Galactic reddening law with  $R_V = 3.1$  (Cardelli et al. 1989). However, the user can choose other extinction laws and  $R_V$  values, if desired, given the code’s open-source nature. If the input measurements are already de-reddened, the user can easily disable the reddening correction. If either  $H\alpha$  or  $H\beta$  are not provided, the reddening correction cannot be implemented. The user is notified with a warning message and has the option to proceed with the calculations with uncorrected line fluxes.

We obtain  $12 + \log_{10}(\frac{O}{H})$  values and their uncertainties in the metallicity scales listed below, and the user can

<sup>6</sup> [https://github.com/nyusngroup/MC\\_Metallicity/blob/master/README.md](https://github.com/nyusngroup/MC_Metallicity/blob/master/README.md)

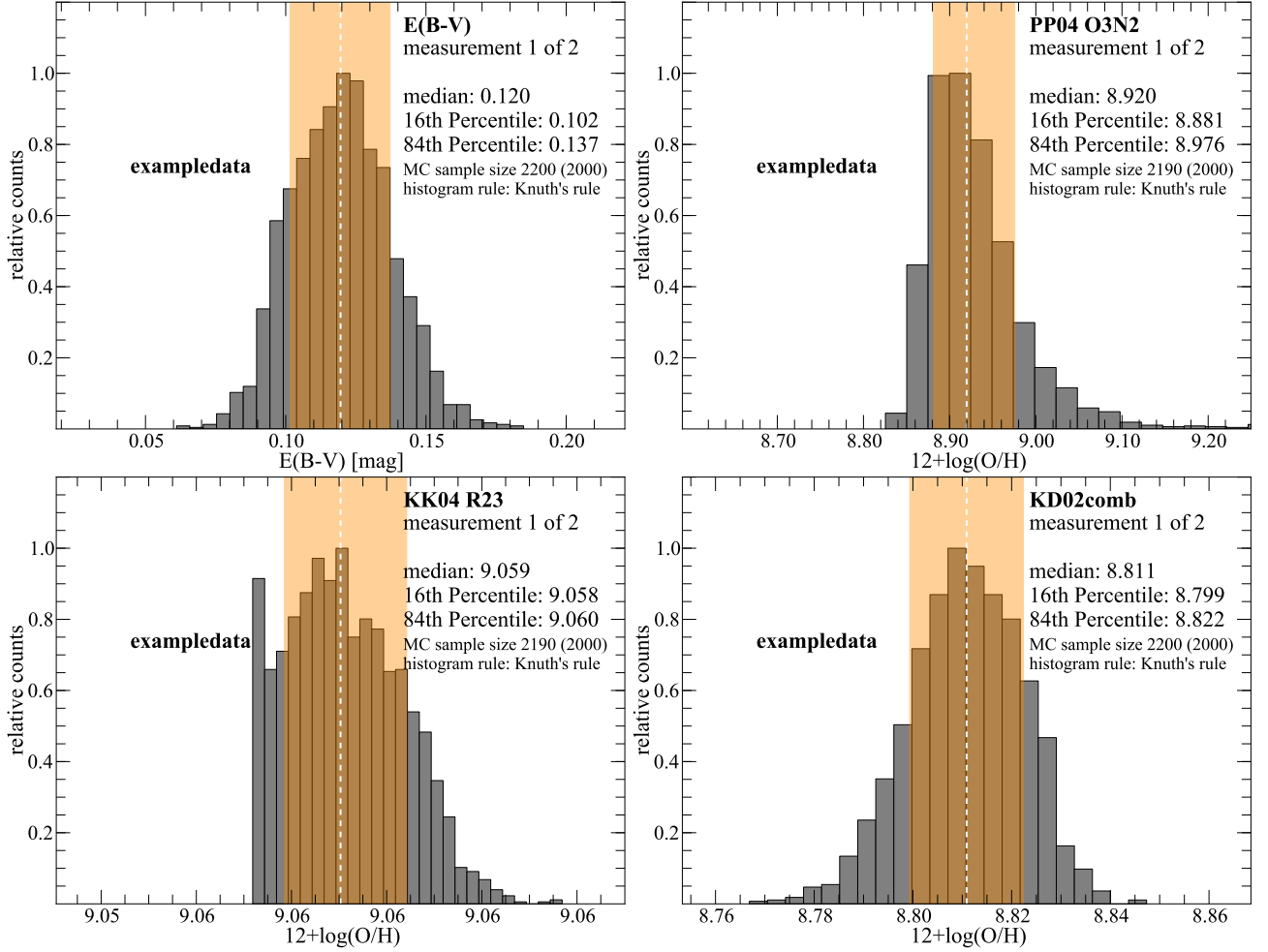


FIG. 1.— Metallicity and reddening  $E(B-V)$  parameter distributions based on the example data shown in Table 1: emission line data of the HII regions at the position of SN 2008D, published in Modjaz et al. (2011). The distributions are generated from  $N=2,000$  samples. The median values are shown with the dashed white lines, while the confidence region, between the 16<sup>th</sup> and the 84<sup>th</sup> percentile is shaded (orange). We show the metallicity scales from Pettini & Pagel (2004), using OIII and NII (PP04 O3N2), the  $R_{23}$ -based scale described in Kobulnicky & Kewley (2004) (KK04\_ $R_{23}$ ), and the combined scale of Kewley & Dopita (2002) (KD02\_comb), updated as described in Kewley & Ellison (2008). Similar plots are outputted by the code for all computed metallicity scales. Each plot indicates: the scale, the sequential number of measurement in input (which corresponds to a line of the input ASCII files), the median, 16<sup>th</sup> and 84<sup>th</sup> percentile values, the sample size (which if by default 10% larger than the requested  $N$  value, but can be smaller if some simulations lead to invalid metallicities), and, finally, the method used to choose the bin size for the histogram (Kuth’s rule in this case, see Section 2.2.1).

choose which of the following scales to calculate, which have been implemented as prescribed in Kewley & Ellison (2008), except where noted.

- **M91** (McGaugh 1991)
- **Z94** (Zaritsky et al. 1994) which is valid for the upper branch of  $R_{23}$  only, and we conservatively constrain it to  $\log(R_{23}) < 0.9$ , i.e., the range that is covered by the photoionization model grids.
- **D02** (Denicoló et al. 2002) for which we include, in addition to the uncertainties in the measurements, the uncertainty on the fit parameters published in D02.
- **PP04**: (Pettini & Pagel 2004) 2 scales, PP04\_N2, based on the  $[\text{NII}]/\text{H}\alpha$  ratio, and PP04\_O3N2, based on  $(\frac{[\text{OIII}]}{\text{H}\beta} / \frac{[\text{NII}]}{\text{H}\alpha})$ .
- **P05** (Pilyugin & Thuan 2005)
- **KD02 & KK04**: 4 scales: KD02\_N2O2, which uses the  $[\text{NII}]/[\text{OII}]$  ratio (Kewley & Dopita 2002), KK04\_N2H $\alpha$  which uses the  $[\text{NII}]/\text{H}\alpha$  ratio (Kobulnicky & Kewley 2004), KK04\_ $R_{23}$  (Kewley & Ellison 2008, appendix A2.2), which is based on the value of  $R_{23}$ , and a combined method, KD02\_comb that chooses the optimal method given the input line fluxes and is implemented as described in Appendix 2.3 of Kewley & Ellison (2008).
- **M08** (Maiolino et al. 2008): This scale is a combination of the KD02 photoionization models at high metallicities and  $T_e$  based metallicities at low metallicities. Our default outputs their strong line diagnostic that is based on  $R_{23}$ , and the diagnostics based on O3O2 and  $[\text{NII}]/\text{H}\alpha$ , since the metallicity estimates from  $[\text{NII}]/\text{H}\alpha$  or O3O2 are necessary to resolve the degeneracy in the double valued  $R_{23}$  metallicity. The other scales (based on  $[\text{OII}]/\text{H}\beta$ ,  $[\text{OIII}]/\text{H}\beta$ , and  $[\text{OIII}]/[\text{NII}]$ ) can be outputted upon explicit user request, via the command

line options.

- **P10** (Pilyugin et al. 2010): This is the so-called ONS diagnostic (involving the [OII], [OIII], [NII] and [SII] lines) and is calibrated with HII regions that have  $T_e$  based metallicities.
- **M13** (Marino et al. 2013): Two scales: M13\_N2, which is a linear fit to the  $\frac{[\text{NII}]6584}{\text{H}\alpha}$ , and M13\_O3N2. This is an updated calibration of the PP04 O3N2 method, based on a large number of  $T_e$ -based metallicity measurements, including those from the CALIFA survey (almost three times larger than the sample used in PP04). This method derives a significantly shallower slope between the O3N2 index and the oxygen abundance than the PP04 calibration did.
- **D13**: (Dopita et al. 2013). The photoionization models used in KD02 and in KK04 have been updated in Dopita et al. (2013) by including new atomic data within a modified photoionization code and by assuming a  $\kappa$  distribution for the energy of the electrons in the HII regions, rather than the simple Maxwell-Boltzmann distribution assumed in prior works. This distribution is more realistic, as observed in Solar System astrophysical plasma (Nicholls et al. 2012). If the user has installed their publicly available `pyqz`<sup>7</sup> Python module, [NII], [SII], [OIII],  $\text{H}\alpha$ , and  $\text{H}\beta$  lines are fed to the `pyqz` module, which produces up to 8 emission line ratio diagnostics for  $12 + \log_{10}(\frac{\text{O}}{\text{H}})$ , each using two of the line ratios [NII]/[SII], [NII]/ $\text{H}\alpha$ , [OIII]/[SII], and [OIII]/ $\text{H}\beta$ . Our code sets the  $\kappa$  parameter to 20, which is the value that Dopita et al. (2013) found best resolves the inconsistencies between oxygen abundance values derived from the strong-line methods and from the “direct”  $T_e$  method<sup>8</sup>.
- **DP00, P01 & C01 (only upon request)**: (Díaz & Pérez-Montero 2000, Pilyugin 2001, Charlot & Longhetti 2001). DP00 is based on  $S_{23}$ , and is the only  $S_{23}$  scale implemented in the current version of the code. P01 is superseded by P05. C01 produces a diagnostic based on  $R_{23}$ ,  $\text{C01\_}R_{23}$ , and one based on [NII]/[SII],  $\text{C01\_N2S2}$ ;  $\text{C01\_}R_{23}$  was used to calculate the combined scale described in KK04, and included in the IDLKD02 code, but it is no longer used in the new combined scale KD02\_comb, which supersedes the old one. Thus, P01 and C01 are only included for historical reasons. These three scales are *not* part of the default output of our code, however they are still available upon explicit user request, via command line input.

The following diagnostics were calculated in the original IDLKD02 code, and are discussed in detail in Kewley & Dopita (2002) and Kewley & Ellison (2008): C01, P01, M91, Z94, D02, PP04, P05, KD02, KK04 and

KD02\_comb. We refer the reader to those papers for further details.

DP00 is the only diagnostic that relies on sulfur ratios:  $S_{23} = ([\text{SII}]6717 + [\text{SIII}]9069)/\text{H}\beta$ . The shortcomings of  $S_{23}$  as a tool to measure abundances are discussed in Kewley & Dopita (2002). The DP00 diagnostic we implement is corrected with the addition of a term  $\propto (c + S_{23}^3)^{-1}$  as suggested by Kewley & Dopita (2002), which corrects the tendency of the scale to systematically underestimate the abundance, with a discrepancy growing larger at higher metallicity. However, we point out that the scatter in metallicity derived from this diagnostic compared to others remain high.

The distributions of  $E(\text{B-V})$  and  $\log(R_{23})$  values are also part of the default output. While certain parameters, such as the ionization parameter  $q$  and the electron density (using the SIII lines) are computed as long as the necessary lines are provided, they are not outputted in the current version of our code – however, the reader can easily modify the code to suite their needs, given it is an open-source code.

## 2.2. Computing Statistical Uncertainties

The novel aspect of our work is that for every set of input line measurements we introduce a Monte Carlo (MC) resampling method to obtain iterations via random sampling within the measurement errors, and thus we obtain a robust result for error estimation (e.g., Efron 1979; Hastie et al. 2009; Andrae 2010).

Given a data set with error bars from which certain parameters are estimated, Monte Carlo resampling generates synthetic data samples drawing from a given distribution. Here we draw synthetic data from a Gaussian distribution centered on each measured line flux value, with standard deviation corresponding to the measurement error. The implicit assumption is made, of course, that the line flux error is Gaussian distributed in nature<sup>9</sup>. For each metallicity scale, and for each of the  $N$  values chosen randomly within the relevant emission line distributions we run the calculations that computes the metallicity. By generating synthetic data, this method effectively simulates conducting multiple experiments when repeating observations is impractical or impossible, as is the case of the emission line flux data.

The sample size  $N$  is set by the user, and one should expect an appropriate value of  $N$  to be a few 1000s, depending on the metallicity scale chosen and measurement errors (for example  $N = 2,000$  is determined to be sufficient for our example data, as shown in Section 2.3, and we provide a tool to assure the sample size is sufficiently large, which we describe in Section 2.3).

A distribution of parameter estimates for the oxygen abundance is generated, from which the median metallicity and its confidence region are calculated, and the results are binned and visualized in an outputted histogram (see Section 2.2.1). This is done for each scale the user chooses to calculate. The fiftieth (50%) percentile, i.e. the median, is reported as the estimated metallicity value, and the 16th and 84th percentiles of the metallicity estimate distribution, as its confidence region. The user

<sup>7</sup> <https://datacommons.anu.edu.au:8443/DataCommons/item/anudc:5037>

<sup>8</sup> The user can modify the value of  $\kappa$  by editing the code, if they wish.

<sup>9</sup> Users may wish to provide their own probability distribution for the emission line uncertainties, and modify the code to suite their needs.

can choose to also output the full metallicity distribution as a binary<sup>10</sup> or an ASCII files.

This MC resampling approach takes into account the impact of the uncertain reddening (due to the uncertainties in the measurement of the  $H\alpha$  and  $H\beta$  fluxes), when the option for de-reddened metallicities is chosen. For each synthetic set of measurements a new reddening value is calculated based on the resampled  $H\alpha$  and  $H\beta$  fluxes, and used to compute the de-reddened metallicity value, thus the derived distribution of metallicity values naturally takes into account the uncertainty in reddening. As part of the code output, the median, confidence intervals are provided, as well as a distribution histogram for  $E(B-V)$  (first panel in figure 1). If either  $H\alpha$  or  $H\beta$  fluxes are not provided, no reddening correction can be applied. The computed metallicity will not be reddening-corrected and the  $E(B-V)$  output will be set to zero.

Figure 1 shows metallicity estimate distributions for three representative scales (PP04\_O3N2, KK04\_ $R_{23}$ , and KD02.comb), and for the reddening parameter  $E(B-V)$  - similar plots that are out-putted by our code for all scales calculated, and for  $R_{23}$ . Although the input line flux distributions are Gaussian, the metallicity distributions are not, for two reasons: first, since the metallicities are computed based on log values of line flux ratios, symmetric error bars in linear space will translate into asymmetric error bars in log space; and second, some metallicity scale computations are non-linear, and sometimes bimodal, especially those that use  $R_{23}$  and  $S_{23}$ , as shown in Figure 2, since the upper or lower branch metallicity value has to be chosen to break the degeneracy for each synthetic measurement.

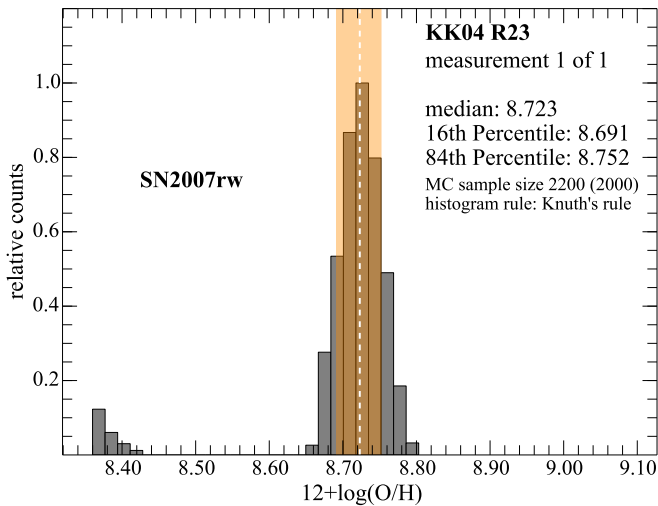


FIG. 2.— Metallicity distribution according to the KK04\_ $R_{23}$  scale for the site of SN 2007rw, as measured in Modjaz et al. (2011). The  $R_{23}$  based scales are double valued: for the same  $R_{23}$  value there are two metallicity solutions. The KK04\_ $R_{23}$  scale uses the ionization parameter  $q$  in an iterative fashion to determine the metallicity by selecting the upper or the lower branch value. In some cases, when the errors in the measurements are large, or for particular flux ratios, the solution may oscillate between the upper and lower branch in different realizations, giving rise to a bimodal metallicity distribution. These cases are easily identifiable by looking at the visual tools the code generates, such as this histogram.

<sup>10</sup> using the *pickle* Python module

Since the metallicity distributions are not Gaussian, the percentiles we report cannot be expressed in terms of  $\sigma$  values. In determining the confidence region for asymmetric and multi-modal distributions, there are broadly three approaches (e.g., Andrae 2010): choosing a symmetric interval, the shortest interval, or a *central* interval. With the central method we determine the confidence interval by choosing the left and right boundaries such that the region outside the confidence interval on each side contains 16% of the total distribution – in analogy to the one-sigma-interval of a Gaussian distribution. This ensures that the algorithm finds the proper boundaries even for asymmetric, non-Gaussian distributions, and in the case of multiple peaks. In summary, the output for the measured value corresponds to the fiftieth (50%) percentile, while the lower error bar corresponds to the 50<sup>th</sup>-16<sup>th</sup> percentile and the upper error bar corresponds to 84<sup>th</sup>-50<sup>th</sup> of the metallicity estimate distribution.

The distributions for the D02 scale include the uncertainty in the fit parameters: the oxygen abundance in this scale is generated as  $12 + \log_{10}(\frac{O}{H}) = 9.12 \pm 0.05 + (0.73 \pm 0.10) \times [NII]/H\alpha$  (Denicoló et al. 2002). The parameters of the fit are generated as the sum of the nominal parameters (9.12 and 0.73) and a Gaussian distributed random value centered on zero, and within a standard deviation of 0.05 and 0.10, respectively, in the above units. Similarly, the distributions for the M13 scales include the uncertainty in the fit parameters as stated in Marino et al. (2013): the oxygen abundance as a function of  $[NII]/H\alpha$  is parametrized as  $12 + \log_{10}(\frac{O}{H}) = 8.743 \pm 0.027 + (0.462 \pm 0.024) \times [NII]/H\alpha$ , and as a function of  $\frac{[OIII]/H\beta}{[NII]/H\alpha}$  as  $12 + \log_{10}(\frac{O}{H}) = 8.533 \pm 0.012 + (0.214 \pm 0.012) \times \frac{[OIII]/H\beta}{[NII]/H\alpha}$ .

We note that our code does not output the *systematic* uncertainty of each scale, which, for example, is  $\sim 0.15$  dex for KD02. However, if all metallicity measurements are in the *same* scale and only *relative* comparisons are made, as recommended by a number of authors, then the systematic error does not have any impact (by definition!). Nevertheless, it is then necessary to obtain statistical error bars, as computed in the code we are releasing here, for meaningful relative comparisons.

### 2.2.1. Visual diagnostics

In order for the user to check the validity of a measurement, and to better understand the distribution of metallicity values, we provide two visualizations: for each set of input line fluxes, we generate a histogram of the output distribution in all metallicity scales calculated (Figures 1, 2, and 3), and for each set of input line fluxes we generate a *box-and-whiskers* plot (hereafter *boxplot*, for short) summarizing the result of all scales calculated (Figure 4). All the plots generated by our code are created with Python *matplotlib* (Hunter 2007).

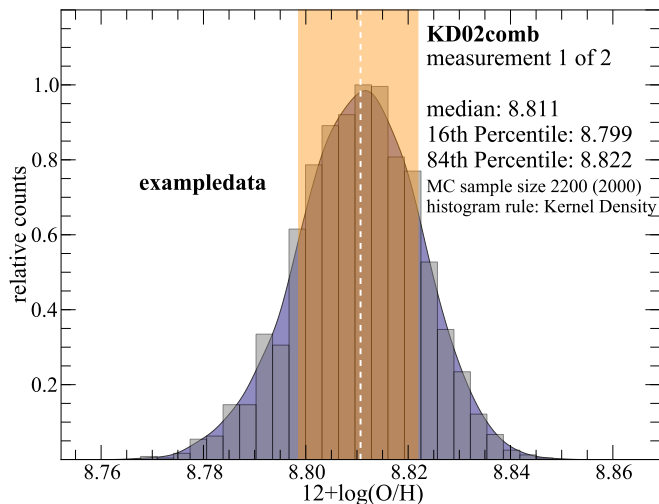
Choosing the binning size for a histogram is not a trivial task. Hogg (2008) describes various data analysis recipes for selecting a histogram bin size. Too many bins will result in an “over-fit” histogram, while too few bins may miss features of the distribution. By default, we use *Knuth’s Method* to choose the number of bins  $N_{bins}$  for each histogram. Knuth’s method optimizes a Bayesian



TABLE 1  
EXAMPLE DATA AND THEIR UNCERTAINTIES BASED ON DATA IN MODJAZ ET AL. (2011)

site	[OII]3727	H $\beta$	[OIII]4959	[OIII]5007	[OI]6300	H $\alpha$	[NII]6584	[SII]6717	[SII]6731
08D	1.842 (.053)	0.958 (.032)	NaN	0.302 (.029)	0.127 (.021)	4.746 (.026)	1.642 (.022)	0.941 (.021)	0.543 (.019)
06fo	2.875 (.101)	1.251 (.044)	0.168 (.028)	0.064 (.025)	NaN	4.026 (.069)	0.781 (.033)	0.821 (.035)	0.573 (.032)

fitness function across fixed-width bins (Knuth 2006). Additionally, however, we enable a number of binning options from which the user can choose, including: the square root of the number of points,  $N_{\text{bins}} = \sqrt{N}$ , *Rice rule* ( $N_{\text{bins}} = 2\sqrt[3]{N}$ , e.g., Hastie et al. 2009), *Doane's formula* ( $N_{\text{bins}} = 1 + \log_2 N + \log_2 \left(1 + \text{Kurt} \sqrt{N/6}\right)$ ), where Kurt is the third moment of the distribution, Doane 1976<sup>11</sup>), and the full Bayesian solution, known as Bayesian Blocks, which optimizes a fitness function across an arbitrary configuration of bins, such that the bins are of variable size (Scargle et al. 2013). The implementation of the latter method requires the `astroML` Python package to be installed on the user's system (Vanderplas et al. 2012<sup>12</sup>). If the `astroML` package is not found, the code will default to Knuth's Rule. As mentioned, Knuth's method implies an optimization. In cases in which the convergence takes too long (or if the number of bins after the minimization is  $N_{\text{bins}}/\sqrt{N} > 5$  or  $N_{\text{bins}}/\sqrt{N} < 1/3$ ) the code will revert to Rice rule. Some methods may be computationally prohibitive with a very large sample size, or very little computational power, such as the Bayesian methods that rely on optimization (Knuth's method and the Bayesian Block method) in which case it may be convenient for a user to choose Doane's formula, Rice rule, or even the square root of the number of samples.



<sup>11</sup> Doane (1976) attempted to address the issue of finding the proper number of bins for the histogram of a skewed distribution. Several versions of the so-called Doane's formula exist in the literature. Our formula can be found, for example, in Bonate 2011

<sup>12</sup> <https://github.com/astroML/astroML>

FIG. 3.— The Kernel Density for the distribution of values for the KD02comb scale, as described in Kewley & Ellison (2008), for the first measurement of our example data (Table 1) is plotted. The Kernel density is displayed as a blue shaded region, and it is calculated via *KD Tree* with a gaussian kernel with bandwidth given by Silverman's rule as described in Section 2.2.1, and normalized as described in the same Section. The histogram of the distribution with min size chosen according to Knuth's rule is over-plotted (gray bins) and the median and confidence intervals are shown as described in Figure 1.

Lastly, the user can generate and visualize the metallicity distribution *Kernel Density* if the `sklearn` package is installed. Kernel Density Estimation (KDE) alleviates the problem of choosing the bin size, at the cost of specifying a convolution kernel. The Kernel Density of the distribution is here calculated via *KD Tree* with a gaussian kernel, as explained in the `sklearn` package documentation<sup>13</sup>. The bandwidth of the kernel is chosen accordingly to *Silverman's rule* (Silverman 1986)<sup>14</sup>. The results will then show both a histogram, with  $N_{\text{bins}}$  chosen via Knuth's method, as well as the distribution Kernel Density, as shown in Figure 3. The KDE is saved as binary (`pickle`) python object (an `sklearn KernelDensity` object), so that it can be recovered outside the code and used as a probability distribution for further inference.

The histograms are always normalized so that the highest bin value is 1, and the Kernel Density is normalized to contain the same area as the overplotted histogram.

At each run the code also generated a boxplot (Figure 4). The boxplot summarizes the result from each scale the user chooses to calculate. For each scale the median of the  $12 + \log_{10}(\frac{O}{H})$  distribution is plotted as a black horizontal line. The height of the corresponding box represents the 25<sup>th</sup> percentile of the  $12 + \log_{10}(\frac{O}{H})$  distribution, known as the *interquartile range* (IQR). The bars represent the maximum and minimum value of the distribution, excluding outliers. The outliers are plotted as circles, and are defined as all data points farther than  $1.5 \times \text{IQR}$  from the 25<sup>th</sup> percentile (i.e. from either end of the box).

The Solar oxygen abundance is indicated in this plot for comparison: a gray box shows a range of estimated values for Solar oxygen abundances, from  $12 + \log_{10}(\frac{O}{H}) = 8.69$  (Asplund et al. 2009) to  $12 + \log_{10}(\frac{O}{H}) = 8.76$  (Caffau et al. 2011). Notice that only the diagnostics requested by the user have a slot in the plot (in the example in Figure 3 the computed scales are Z94, M91, the PP04, the M13, and the KD02 scales). However these slot exists on the plot whether the diagnostic

<sup>13</sup> <http://scikit-learn.org/stable/modules/density.html>

<sup>14</sup> By Silverman's rule the bandwidth for the KDE kernel is set to  $w = 1.06\sigma N^{-0.2}$ , where  $\sigma$  is the sample standard deviation. Although the bandwidth chosen accordingly to Silverman's rule is only optimal in the case of a Gaussian basis and a Gaussian distribution, and our distributions are explicitly *not* Gaussian, as explained earlier, this kernel parametrization generally provides good results for our metallicity distributions.

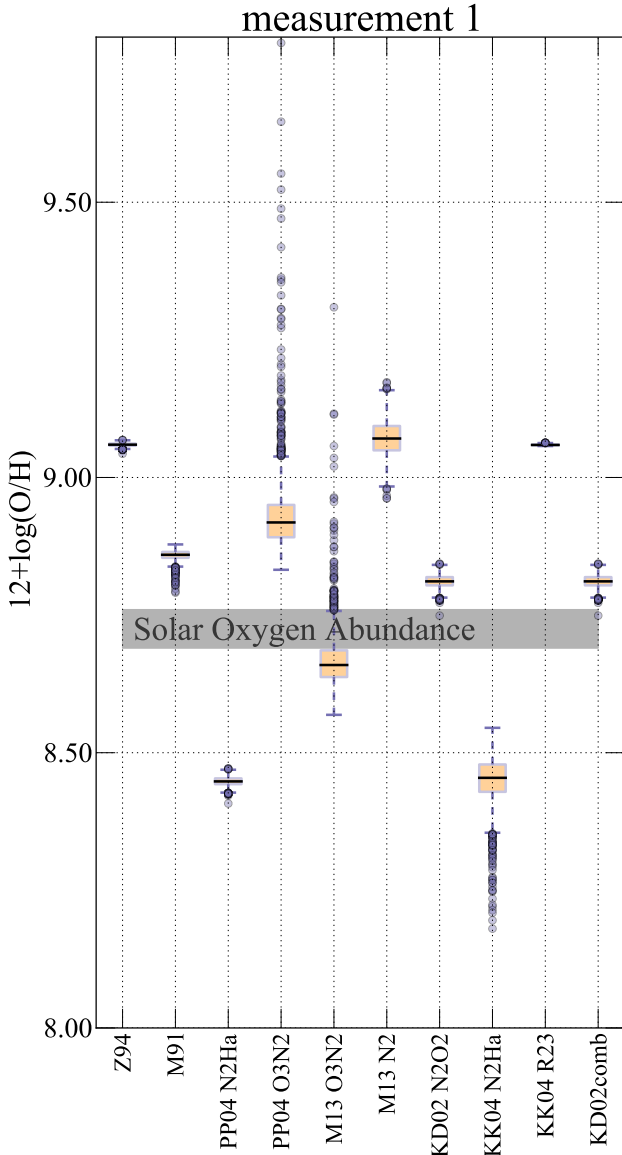


FIG. 4.— A box-and-whiskers plot shows the comparison of the results of 11 metallicity calculations, corresponding to 5 scales as listed above, calculated from the same set of measured lines (Table 1, host galaxy of SN 2008D). For each scale, the median of the resulting distribution is plotted as a horizontal line, the inter quartile range (IQR) is represented as an orange box, and the bars, joined to each end of the box by a dashed line, represent the minimum and maximum of the distribution *excluding outliers*, where outliers are defined as any point farther than  $1.5 \times \text{IQR}$  from the edges for the IQR.

can be produced or not, i.e. if the set of input lines does not allow a requested scale to be calculated an empty column will be generated in this plot in correspondence of said metallicity scale.

With this plot the user can immediately check for consistency or scatter in the metallicities derived by the requested scales (e.g., there are well-known systematic offsets between different methods, where the  $T_e$ -based, empirical methods usually give lower values than the photoionization based ones), and coarsely assess the shape of the distribution in each scale (e.g. strong asymmetry or bimodality would show up as an asymmetric box, or a very asymmetric distribution of outliers).

### 2.3. Visual diagnostic to assure completeness of the MC simulation

The user chooses  $N$ , the number of samples to be generated. Note that the final sample size is 10% larger than the user chosen parameter  $N$ . The sample size is in fact automatically increased by 10% at the beginning of the run, in order to assure that even if during the calculations some of the output metallicities were to result in non-valid values ( $NaN$ 's or infinities, for example, if division by very small numbers is required) the actual sample size is at least as large as the user intended. With reasonable input parameters, the code rarely produces non-valid values, and if the size of the valid output distribution were in the end smaller than the *requested* value  $N$ , i.e. if the number of non-valid outputs is larger than 10% of  $N$ , the user should in fact worry about the set of input line fluxes used.

Of course the reliability of the metallicity estimates depends crucially on the sample size being sufficiently large to properly characterize the distribution of metallicity values. It is however not trivial to decide when  $N$  is sufficiently large. As soon as  $N$  is large enough, and the distribution is well characterized, adding more synthetic data will not change its shape. Consider a cumulative distribution for a metallicity scale, D02, for example, which has noise from the measurement errors, as well as from the error in the fit parameters, or KD02, which uses a non linear combination of the input line flux values. We plot the cumulative distribution for four randomly selected subsamples of the data points in the distribution, of size 10%, 25%, 50%, 75% of  $N$ , as well as for the entire distribution. If 1/10th of the points were a sufficiently large Monte Carlo sample, the cumulative distributions would typically appear smooth and, most importantly, they would overlap. In Figure 5 we show these cumulative plots for a distribution generated with  $N = 200$ ,  $N = 2,000$ , and  $N = 20,000$ , for D02 and KD02comb.

The cumulative distributions for the subsets of the  $N = 200$  simulation (left panel in Figure 5) are different, and noisy, indicating that 200 data points is not a sufficiently large dataset. Conversely, at  $N = 20,000$  (right panel) the distributions are indistinguishable, indicating that  $N = 2,000$  (the smaller subset in this panel) is large enough. The  $N = 2,000$  cumulative distributions rapidly converge, as the subsample size increases, indicating that a value of  $N$  between 200 and 2,000 will characterize the distribution appropriately for our example data set. In the light of this we choose, conservatively,  $N = 2,000$  to run our simulations for this dataset.

We remind the user that the appropriate number of  $N$  samples will depend on the diagnostic, and on the input data. The errors on the measurements, and the calculations performed to derive the metallicity from line ratios, which are for many scales non-linear, will determine the shape of the distribution, and thus the number of data points that are needed to fully characterize it.

### 2.4. Performance and Benchmarks

We ran a bench mark calculation on a MacBook Air with a dual-core Intel Core i7 (1.7 GHz) and 8GB of 1600 MHz DDR3 memory. The dataset we used for the calculation includes 2 sets of line measurements (measure-



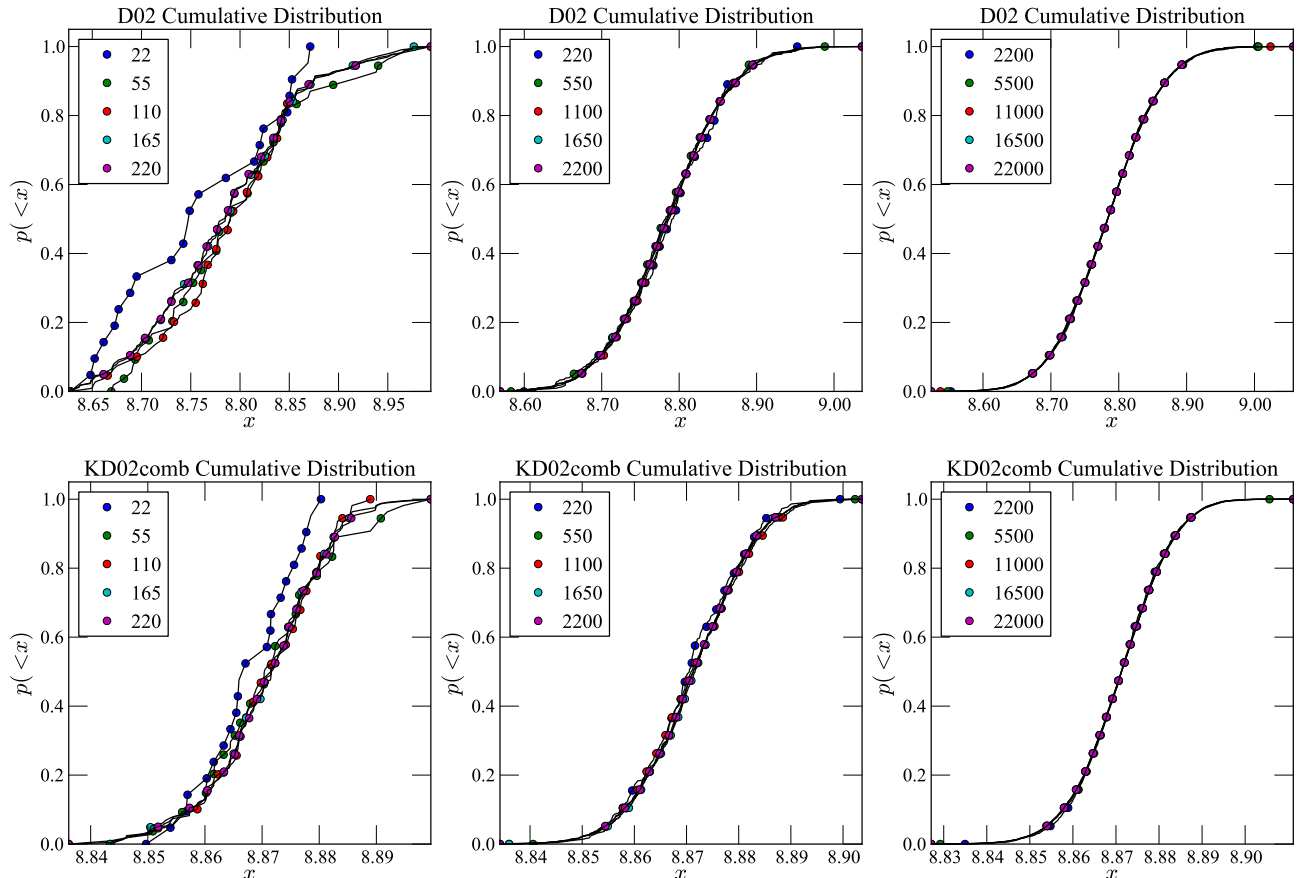


FIG. 5.— Cumulative plots of the distribution of metallicity values for the D02 (Denicoló et al. 2002) and KD02 scale (Kewley & Dopita 2002, updated by Kewley & Ellison 2008), chosen here just as examples, where  $x$  refers to  $12 + \log_{10}(\frac{O}{H})$ . The input used is example data 1, the emission line values for the host galaxy of SN 2008D from Modjaz et al. (2011). In each plot the cumulative distribution is shown for a randomly chosen subsample of 10% of the data in the sample, as well as a subsample of 25%, 50%, and 75%, of the data, and for all data in the distribution. In the left plots the distributions are generated from an  $N = 200$  samples, in the center top and bottom plots from an  $N = 2,000$  and in the right-most column from a  $N = 20,000$  sample. The increasing overlap of the distributions informs us about completeness. In the left plots the distributions do not fully overlap, indicating that completeness is not achieved with the  $N = 200$  sample. On the other end, since all subsamples are indistinguishable in the top and bottom plots on the right, which are generated with  $N = 20,000$  samples, we conclude that completeness is already achieved at 10% of the  $N = 20,000$  sample (the smallest sample in the plots on the right) for our example data, i.e., at  $N = 2,000$  which we use for the rest of this work.

ment 1, the host galaxy of SN 2008D, and measurement 2, the host galaxy of SN 2006fo). The flux values and their associated errors are shown in Table 1. The code performs simple algebraic operations on large data arrays. It is vectorized along the sample-size dimension, so that the code loops over the smaller dimension corresponding to the number of sets of line measurements in input (2 in our example data), while operations are generally performed on the  $N$ -dimensional vectors storing the synthetic measurements, and the  $N$ -dimensional variables hence derived (certain scales, such as the KD02 and KK04 ones, derive the ionization parameter in an iterative fashion requiring further loops).

With full graphic output (all histograms are being plotted) and performing all default metallicity calculations, except those of D13 *pyqz*, for our example data sets and a sample size  $N = 2,000$ , the time required by the code is  $\sim 14.5$  seconds (wall clock time, and less than half a second longer in total CPU time, recall that the machine we tested on is dual-core), and less than 0.3 seconds of CPU time spent in the kernel within the process. Including the scales of D13 *pyqz* for the same datasets, the run time becomes  $\sim 29$  seconds of clock time (and actual

CPU time). The code time was tested on sets of 1, 2, 5, 10, 25, 50 and 100 identical measurements (copies of the emission line fluxes of the SN 2008D host galaxy in our example dataset) and the clock time is found to scale roughly quadratically with the number of measurements in input, but with a small quadratic coefficient  $\sim 0.07$ , and with linear coefficient  $\sim 7$ . This means effectively that the 100 measurements sample will take 200 times longer than the 1 measurement sample (of course with dependence on which lines are available in each measurement). For the CPU time spent in the kernel we find a roughly linear relation with the number of input measurements, with a very shallow dependency of  $\sim 0.05$ . This is the actual computational time in calculations performed in the code: most clock time is in fact spent in root finding, input-output, and plotting activities.

The time spent on plotting functions, which includes the calculation of the appropriate number of bins for each scale, is substantial: 1.67 seconds per distribution on average with 14 calls for this data set (one for each metallicity scale,  $E(B-V)$ , and  $R_{23}$ ). We summarize the run time spent on each metallicity scale and memory usage for the host galaxy of SN 2008D, the first measurement of

the sample dataset, in Figure 6. While the CPU usage is modest, the memory usage can be large, depending on the size of the sample. The memory usage ramps up quickly, as soon as the  $N$  line samples are created, and remains fairly constant throughout the run. Figure 6 shows the memory used by each function in the code as a function of time for a single set of input lines of our example data (host galaxy of SN 2008D). The scales that take longer time require the root finding (e.g. M08), or optimizations which are done iteratively (e.g. KD02\_N2O2).

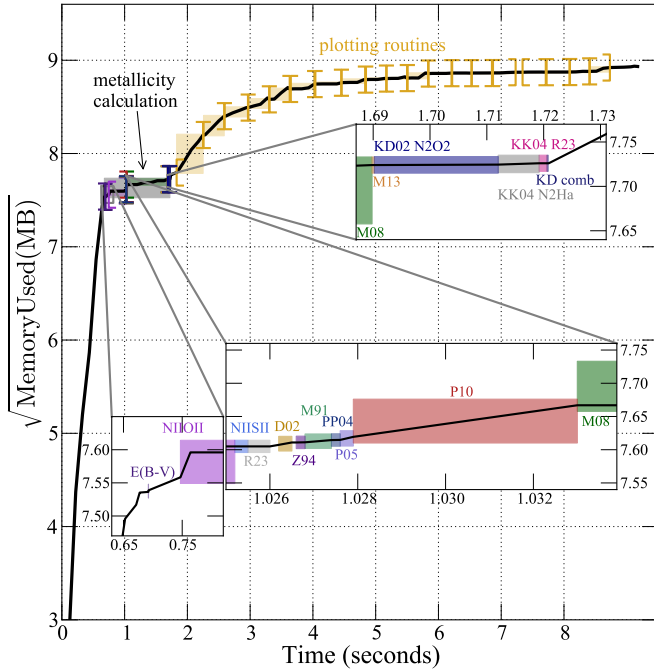


FIG. 6.— Memory usage: we plot the square root of the memory usage in Megabytes as a function of time for running our code (using  $N=2,000$  and all default metallicity scales except the D13 *pyqz* ones) on a single set of measured emission lines (Table 1, host galaxy of SN 2008D). The square root is plotted, instead of the natural value, to enhance visibility. Three inserts show zoom-ins of the region where most of the metallicity scales are calculated, since the run time of the code is dominated after all scales are computed, by plotting routines, including the calculation of the bin size with Knuth’s rule. Each function call is represented by an opening and closing bracket in the main plot, and by a shaded rectangle in the zoom-in inserts. The calculation of NII0II, which requires 0.26 seconds (update in the end), is split between two inserts, as well as the calculation of the M08 scales (three scales) which require 0.78 seconds. Altogether, the calculation of metallicity in all scales, except those in the D13 *pyqz* scales, takes 1.15 seconds. For technical details about the D13 *pyqz* scales performance we refer the reader to the *pyqz* package. **FIX WITH FINAL VERSION OF PLOT**

The code is trivially parallelizable, as each measurement can be computed on a different core. Multithreading is enabled via the `multiprocessing` Python module, and when enabled it uses up the  $n - 1$  codes, where  $n$  is the total number of available cores to the user (or less, if a smaller maximum number of cores is set by the user via the global variable `MAXPROCESSES`).

## 2.5. Availability

The source code is published under MIT-license<sup>15</sup> on GitHub<sup>16</sup>. At this time the code is released under DOI NUMBER HERE as version 1.0: NAME HERE. Project details, step-wise tutorials, and further information can be found in the module README<sup>17</sup>. Development is done in Linux and OS X. The package requires standard python packages, such as `numpy`, `scipy`, `pylab`, and additional features are enabled if the packages `astroml`, `skitlearn`, and `pyqz` are installed, but these packages are *not* required. Contact the authors to be included in a mailing list and be notified about critical changes.

## 3. COMPARISON TO PRIOR UNCERTAINTY COMPUTATION AND OTHER WORKS

A previous method for determining the uncertainty in the oxygen abundance (as used in Modjaz et al. 2008; Kewley et al. 2010; Rupke et al. 2010; Modjaz et al. 2011) was an *analytic* approach of propagating the emission-line flux uncertainties: it found the maximum and minimum abundances via minimizing and maximizing the line ratios by adding/subtracting to the measured line flux values their uncertainties. For comparison we computed the metallicities and their errors in both ways (analytically and using our current MC resampling method) for 3 representative scales. We plot our results and the residuals in Figure 7, which shows a number of important points: i) The metallicity reported as the 50<sup>th</sup> percentile of the metallicity parameter distribution from the MC resampling method is completely consistent with the analytically derived metallicity (i.e., **DAVID, please fill in** with XXX, XXX and XXX dex for the standard deviation of the residuals for the KD02-comb, M91 and PP04-O3N2 scales, respectively) - well within the respective error bars - and thus, the prior published results still stand (unsurprisingly, since our code, aside for the calculation of the confidence interval, uses the same algorithms developed for IDLKD02). ii) The MC resampling method has smaller error bars than the analytic method, especially for the scales of M91 and KD02. This is easily understandable, since the analytic method assumes the worst-case-scenarios, as it basically yields 2 metallicity parameter draws (the “minimum” and “maximum”) which are in the tail of the full metallicity probability distribution. The MC resampling method is a more correct method as it empirically characterizes the full parameter estimation distribution.

### 3.1. Comparison with other works

The field of SN host metallicity studies has been rapidly developing as these kinds of studies may be crucial avenues for constraining the progenitor systems of different kinds of explosions - however, a few of the works in this field do not compute errors and others do not show how they compute their statistical errors (e.g., Anderson et al. 2010; Leloudas et al. 2011; Sanders et al. 2012; Leloudas et al. 2014).

In contrast, the general metallicity field has considered in detail how to estimate the uncertainties in measured

<sup>15</sup> [https://github.com/nyusngroup/MC\\_Metallicity/blob/master/LICENSE.txt](https://github.com/nyusngroup/MC_Metallicity/blob/master/LICENSE.txt)

<sup>16</sup> [https://github.com/nyusngroup/MC\\_Metallicity](https://github.com/nyusngroup/MC_Metallicity)

<sup>17</sup> [https://github.com/nyusngroup/MC\\_Metallicity/blob/master/README.md](https://github.com/nyusngroup/MC_Metallicity/blob/master/README.md)

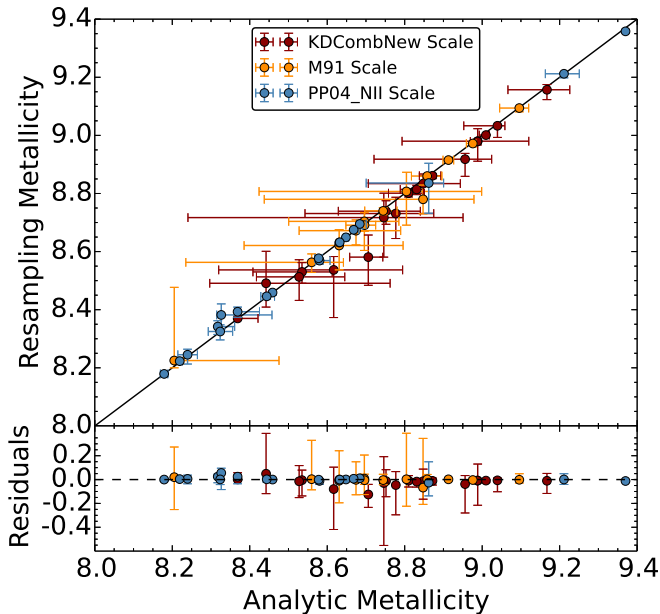


FIG. 7.— Comparison of metallicity estimation between the analytic method and our Monte Carlo resampling method (top) and their residuals (bottom) for three different metallicity scales. Flux measurements come from 19 galaxies previously measured in Modjaz et al. (2011). To add asymmetric errors in quadrature we use

$$\text{residual}_{\min} = \sqrt{x_{\max}^2 + y_{\min}^2} \text{ and } \text{residual}_{\max} = \sqrt{x_{\min}^2 + y_{\max}^2}.$$

Note that for the scales of KD02\_comb and M91, the MC resampling error bars are smaller than those of the analytic propagation, which assumes worst-case scenarios, while metallicity reported as the 50th percentile of the metallicity parameter distribution from the MC resampling method is completely consistent with the analytically derived metallicity in all scales.

metallicities- however, none of those codes are open-source and many of them are for specific scales which were chosen by the authors: Moustakas et al. (2010) also use MC resampling to estimate the metallicity uncertainties (in their case using  $N=500$  trials and assuming a Gaussian distribution for metallicity parameter distribution) but only do this for two scales, KK04 and P05. For computing the metallicities of the SDSS star forming galaxies, Tremonti et al. (2004) fit a combination of stellar population synthesis models and photoionization models to the observed strong emission lines [OII],  $H\beta$ , [OIII],  $H\alpha$ , NII and SII and report the median of the metallicity likelihood distribution as the metallicity estimate, with the width of the distribution giving the  $1-\sigma$  (Gaussian) error. However, this constitutes yet a different metallicity scale (the T04 scale).

In the last stages of preparing this manuscript, Blanc et al. (2015) was published. Blanc et al. (2015) employ Bayesian inference for doing something similar to Tremonti et al. (2004) - they use Bayesian inference to derive the joint and marginalized posterior probability density functions for metallicity  $Z$  and ionization parameter  $q$  given a set of observed line fluxes and an input photoionization model. They provide a publicly available IDL implementation of their method named *IZI* (inferring metallicities -  $Z$  - and ionization parameters) on the author's web site.

#### 4. CONCLUSIONS

We have presented the open-source Python code MCZ for the determination of the strong-emission-line estimators of oxygen abundance and its error distribution in a total of 14 scales (11 default and 3 additional upon request), for a total of up to 33 different metallicity measurements. These estimates are based on the original IDL-code in Kewley & Dopita (2002) and expanded to include more recently developed scales, and to generate, via Monte Carlo resampling, a confidence interval for each measurement. In addition our code supplies visualization tools that enable the user to assess the validity of each derived metallicity distribution, and understand when line flux measurements may lead to misleading metallicity estimates, for example in proximity of the demarcation between the upper and lower branch for  $R_{23}$  based methods. Our code outputs the full estimated metallicity distribution (on demand), and its Kernel Density. Our code also offers visualization tools to assess the spread of the oxygen abundance in the different scales. The validity of our metallicity measurements and of their confidence regions of course hinges upon generating probability distributions that properly sample the metallicity distribution, given the input parameters and the specific metallicity calculation algorithm. Thus we develop metrics that allow the user to ascertain that the sample drawn in the Monte Carlo simulation is sufficiently large.

This code is open access and we welcome input and further development from the community. We hope that this open-access code will be helpful for the many different fields where gas-phase metallicities are important, including in the emerging field of SN and GRB host galaxies, where either it is not described how they got uncertainties or no error bars are computed. Given its public-access nature, the users are free to include any new metallicity scales and modify any parts and assumptions (e.g. that the line flux errors are Gaussian distributed). However, we re-iterate that the tool we are providing here should be used responsibly - the user is responsible for understanding the strengths and caveats of the various diagnostics and in which ranges and conditions they can be used.

The Modjaz SNYU group at NYU is supported in parts by the NSF CAREER award AST-1352405 and by NSF award AST-1413260. F. B. Bianco is supported by a *James Arthur Fellowship* at the NYU-Center for Cosmology and Particle Physics and Y. Liu by a *James Arthur Graduate Award*. This code made use of several Python Modules, including *Matplotlib* (Hunter 2007). Some plots are produced with public code DOI:10.5281/zenodo.15419 available at [https://github.com/fedhere/residuals\\_pylab](https://github.com/fedhere/residuals_pylab). This research made use of NASA Astrophysics Data System; the NASA/IPAC Extragalactic Database (NED), which is operated by the Jet Propulsion Laboratory, California Institute of Technology, under contract with the National Aeronautics and Space Administration.

#### REFERENCES

- Anderson, J. P., Covarrubias, R. A., James, P. A., Hamuy, M., & Haberman, S. M. 2010, MNRAS, 407, 2660
- Andrae, R. 2010, ArXiv e-prints, arXiv:1009.2755

- Asplund, M., Grevesse, N., Sauval, A. J., & Scott, P. 2009, *ARA&A*, 47, 481
- Baldwin, J. A., Phillips, M. M., & Terlevich, R. 1981, *PASP*, 93, 5
- Berg, D. A., Croxall, K. V., Skillman, E. D., et al. 2015, *ArXiv e-prints*, arXiv:1501.02270
- Blanc, G. A., Kewley, L., Vogt, F. P. A., & Dopita, M. A. 2015, *ApJ*, 798, 99
- Bonate, P. 2011, *Pharmacokinetic-Pharmacodynamic Modeling and Simulation*, SpringerLink : Bücher (Springer)
- Caffau, E., Ludwig, H.-G., Steffen, M., Freytag, B., & Bonifacio, P. 2011, *Sol. Phys.*, 268, 255
- Cardelli, J. A., Clayton, G. C., & Mathis, J. S. 1989, *ApJ*, 345, 245
- Charlot, S., & Longhetti, M. 2001, *MNRAS*, 323, 887
- Denicoló, G., Terlevich, R., & Terlevich, E. 2002, *MNRAS*, 330, 69
- Díaz, A. I., & Pérez-Montero, E. 2000, *MNRAS*, 312, 130
- Doane, D. P. 1976, *The American Statistician*, 30, 181
- Dopita, M. A., Sutherland, R. S., Nicholls, D. C., Kewley, L. J., & Vogt, F. P. A. 2013, *ApJS*, 208, 10
- Efron, R. 1979, *Ann. Stat.*, 7, 1
- Grevesse, N., Asplund, M., Sauval, A. J., & Scott, P. 2010, *Ap&SS*, 328, 179
- Hastie, T., Tibshirani, R., & Friedman, J. 2009, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (Springer Science+Business Media, New York)
- Hogg, D. W. 2008, *ArXiv e-prints*, arXiv:0807.4820
- Hunter, J. D. 2007, *Computing In Science & Engineering*, 9, 90
- Johnson, J. L., & Li, H. 2012, *ApJ*, 751, 81
- Kauffmann, G., Heckman, T. M., Tremonti, C., et al. 2003, *MNRAS*, 346, 1055
- Kelly, P. L., & Kirshner, R. P. 2012, *ApJ*, 759, 107
- Kewley, L. J., & Dopita, M. A. 2002, *ApJS*, 142, 35
- Kewley, L. J., & Ellison, S. L. 2008, *ApJ*, 681, 1183
- Kewley, L. J., Groves, B., Kauffmann, G., & Heckman, T. 2006, *MNRAS*, 372, 961
- Kewley, L. J., Rupke, D., Zahid, H. J., Geller, M. J., & Barton, E. J. 2010, *ApJ*, 721, L48
- Knuth, K. H. 2006, *ArXiv Physics e-prints*, physics/0605197
- Kobulnicky, H. A., & Kewley, L. J. 2004, *ApJ*, 617, 240
- Leloudas, G., Gallazzi, A., Sollerman, J., et al. 2011, *A&A*, 530, A95
- Leloudas, G., Schulze, S., Kruehler, T., et al. 2014, *ArXiv e-prints*, arXiv:1409.8331
- Levesque, E. M., Berger, E., Kewley, L. J., & Bagley, M. M. 2010, *AJ*, 139, 694
- López-Sánchez, Á. R., Dopita, M. A., Kewley, L. J., et al. 2012, *MNRAS*, 426, 2630
- Lunnan, R., Chornock, R., Berger, E., et al. 2014, *ApJ*, 787, 138
- Maiolino, R., Nagao, T., Grazian, A., et al. 2008, *A&A*, 488, 463
- Marino, R. A., Rosales-Ortega, F. F., Sánchez, S. F., et al. 2013, *A&A*, 559, A114
- McGaugh, S. S. 1991, *ApJ*, 380, 140
- Modjaz, M. 2012, in *IAU Symposium*, Vol. 279, *IAU Symposium*, 207–211
- Modjaz, M., Kewley, L., Bloom, J. S., et al. 2011, *ApJ*, 731, L4
- Modjaz, M., Kewley, L., Kirshner, R. P., et al. 2008, *AJ*, 135, 1136
- Moustakas, J., Kennicutt, Jr., R. C., Tremonti, C. A., et al. 2010, *ApJS*, 190, 233
- Nicholls, D. C., Dopita, M. A., & Sutherland, R. S. 2012, *ApJ*, 752, 148
- Osterbrock, D. E. 1989, *Astrophysics of Gaseous Nebulae and Active Galaxies* (Mill Valley: University Science Books)
- Pagel, B. E. J., Edmunds, M. G., Blackwell, D. E., Chun, M. S., & Smith, G. 1979, *MNRAS*, 189, 95
- Pan, Y.-C., Sullivan, M., Maguire, K., et al. 2014, *MNRAS*, 438, 1391
- Pettini, M., & Pagel, B. E. J. 2004, *MNRAS*, 348, L59
- Pilyugin, L. S. 2001, *A&A*, 369, 594
- Pilyugin, L. S., & Thuan, T. X. 2005, *ApJ*, 631, 231
- Pilyugin, L. S., Vílchez, J. M., & Thuan, T. X. 2010, *ApJ*, 720, 1738
- Rupke, D. S. N., Kewley, L. J., & Chien, L.-H. 2010, *ApJ*, 723, 1255
- Sanders, N. E., Soderberg, A. M., Levesque, E. M., et al. 2012, *ApJ*, 758, 132
- Scargle, J. D., Norris, J. P., Jackson, B., & Chiang, J. 2013, *ApJ*, 764, 167
- Silverman, B. W. 1986, *Density Estimation for Statistics and Data Analysis* (Chapman and Hall)
- Simón-Díaz, S., & Stasińska, G. 2011, *A&A*, 526, A48+
- Stasińska, G. 2002, *ArXiv Astrophysics e-prints*, astro-ph/0207500
- Stasińska, G. 2010, in *IAU Symposium*, Vol. 262, *IAU Symposium*, ed. G. R. Bruzual & S. Charlot, 93–96
- Tremonti, C. A., Heckman, T. M., Kauffmann, G., et al. 2004, *ApJ*, 613, 898
- Vanderplas, J., Connolly, A., Ivezić, Ž., & Gray, A. 2012, in *Conference on Intelligent Data Understanding (CIDU)*, 47–54
- Zaritsky, D., Kennicutt, Jr., R. C., & Huchra, J. P. 1994, *ApJ*, 420, 87