

Agroinfiltration contributes to VP1 recombinant protein degradation

Priyen Pillay,¹ Karl Josef Kunert,¹ Stefan George van Wyk,¹ Matome Eugene Makgopa,¹
Christopher Ashley Cullis,² and Barend Juan Vorster¹

¹Department of Plant and Soil Sciences, Forestry and Agricultural Biotechnology Institute (FABI), University of Pretoria, Private Bag X20, Hillcrest, Pretoria, South Africa

²Department of Biology, Case Western Reserve University, Cleveland, OH 44106, USA

Corresponding author: Juan Vorster

Address: Department of Plant and Soil Sciences, University of Pretoria, Private Bag X20, Hillcrest, Pretoria, South Africa

Telephone: +27 (12) 420 4491

Fax: +27 (12) 420 3960

Email: juan.vorster@up.ac.za

Abstract

There is a growing interest in applying tobacco agroinfiltration for recombinant protein production in a plant based system. However, in such a system, the action of proteases might compromise recombinant protein production. Protease sensitivity of model recombinant foot-and-mouth disease (FMD) virus P1-polyprotein (P1) and VP1 (viral capsid protein 1) as well as *E. coli* glutathione reductase (GOR) were investigated. Recombinant VP1 was more severely degraded when treated with the serine protease trypsin than when treated with the cysteine protease papain. Cathepsin L- and B-like as well as legumain proteolytic activities were elevated in agroinfiltrated tobacco tissues and recombinant VP1 was degraded when incubated with such a protease-containing tobacco extract. *In silico* analysis revealed potential protease cleavage sites within the P1, VP1 and GOR sequences. The interaction

modelling of the single VP1 protein with the proteases papain and trypsin showed greater proximity to proteolytic active sites compared to modelling with the entire P1-polyprotein fusion complex. Several plant transcripts with differential expression were detected 24 hr post-agroinfiltration when the RNA-seq technology was applied to identify changed protease transcripts using the recently available tobacco draft genome. Three candidate genes were identified coding for proteases which included the Responsive-to-Desiccation-21 (RD21) gene and genes for coding vacuolar processing enzymes 1a (*NbVPE1a*) and 1b (*NbVPE1b*). The data demonstrates that the tested recombinant proteins are sensitive to protease action and agroinfiltration induces the expression of potential proteases that can compromise recombinant protein production.

Keywords: Agroinfiltration, recombinant protein production, proteases, cysteine proteases, Foot-and-mouth disease, VP1 protein, tobacco

Introduction

In 2014, an experimental drug, called “ZMapp”, was used to treat two medical workers who had contracted the deadly Ebola virus. The unique characteristic of ZMapp is that its constituents have been produced in *Nicotiana benthamiana* plants.¹ *N. benthamiana* is a model plant species widely used for the transient expression of proteins. Tobacco is sometimes compared to the role that the white mouse has played in mammalian studies.²⁻⁵ The *N. benthamiana* genome sequence has further potential to be useful for gene mining, construct design, and for the assessment of target and non-target gene silencing.² A future prospect is also applying RNA-Seq data to fully annotate the tobacco genome and characterize the transcriptome.² The large leaves of *N. benthamiana* and its susceptibility to a

variety of pathogens have been harnessed as a means to transiently express proteins, using either engineered viruses or *Agrobacterium tumefaciens*.⁶ Due to a lower content of secondary compounds interfering in the protein purification process, *N. benthamiana* has been previously applied as a model plant species for heterologous protein expression.⁷ It has also been included as a tool in platforms for the production of recombinant proteins for comparative analyses.⁸

Due to proteolysis caused by protease action, plant-expressed recombinant proteins can possibly undergo either complete or partial proteolytic degradation.⁹⁻¹¹ Such degradation can ultimately result in proteins with altered biological activity or no protein production at all.^{12,13} The identification of such proteases involved, particularly in *Nicotiana* species, has therefore been the subject of several recent investigations. The majority of protease families, which might compromise recombinant protein production in *Nicotiana* species, belong to the aspartic and cysteine protease (papain-like) families and, to a lesser extent, the serine and metallo-protease families.^{14,15} There is further evidence that such recombinant protein degradation might occur during the extraction process *ex vivo* as a result of proteases being released during the tissue disruption process.¹⁶ However, almost all protease families have also been associated with plant senescence.¹⁷ In *Nicotiana* species, the majority of these proteases are of aspartic or cysteine-type and, to a lesser extent, of serine and metallo-type.¹⁸ However, the *N. benthamiana* leaf contains less protease activity than a *N. tabacum* leaf and is therefore preferred for agroinfiltration.¹⁵ It has been recently reported that agroinfiltration can significantly alter the distribution of cysteine (C1A) and aspartate (A1) protease along the leaf age gradient in *N. benthamiana*. This was further related to the level of proteolysis in whole-cell and apoplast protein extracts.¹⁹ Improvements have been found for various recombinant proteins when protease activity was altered including: bovine serum albumin

(BSA), human serum immunoglobulins G (hIgGs), anti-HIV antibodies (2F5) as well as human protease inhibitor, α 1-antitrypsin.^{13,20,21} However, there is a lack of detailed information on whether induction of plant-derived proteases is among the host responses to agroinfiltration.²² Therefore, the identification of proteases induced by agroinfiltration might be a key step in the improvement of recombinant protein production when applying the agroinfiltration technique.

The purpose of this study was to investigate protease-sensitivity of various model recombinant proteins with the aim to first establish any protease sensitivity and secondly to identify possible proteases expressed as a consequence of the agroinfiltration process. We specifically hypothesized that cysteine proteases are among these expressed proteases following agroinfiltration based on a previous finding in our group that recombinant *E. coli* glutathione reductase (GOR) was more stable in agroinfiltrated tobacco leaves engineered with a rice cysteine protease inhibitor (OC-I).²³ In our study, we determined the inherent vulnerabilities of recombinant model proteins derived from the foot-and-mouth disease virus (FMDV) which are the VP1 and P1-polyprotein (P1) as well as *Escherichia coli* (*E. coli*)-derived glutathione reductase (GOR) proteins towards proteolysis. We also applied protein modelling to investigate how interacting residues within VP1 would interact with a cysteine and serine protease (papain and trypsin) either individually or as part of a P1-polyprotein to obtain more information on VP1 stability against protease action. Finally by applying transcriptomic profiling using RNA-Seq, *N. benthamiana* leaves were screened for the transcription of proteases due to agroinfiltration. We found that the recombinant model proteins used were sensitive to cysteine and serine protease degradation and that expression of several types of proteases, including cysteine proteases, increased due to the agroinfiltration of tobacco leaves.

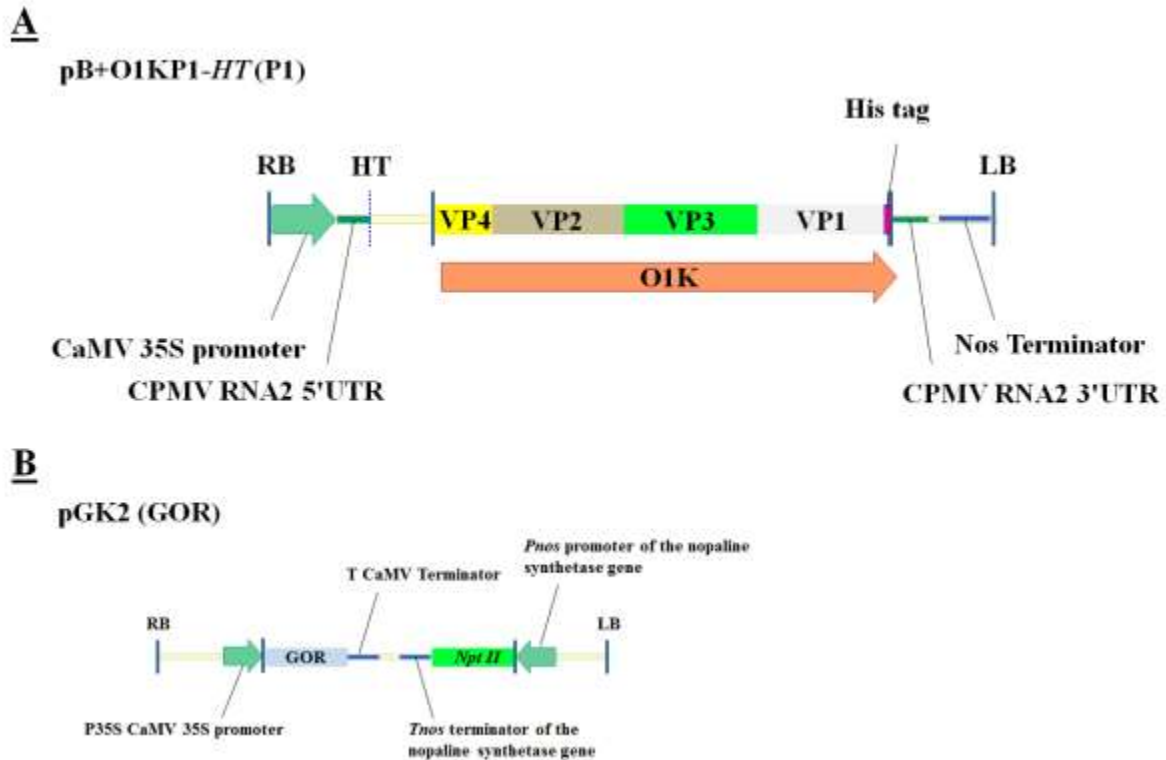


Figure 1 (A) Binary vector pB+O1KP1-HT harbouring the coding sequence O1K under control of a duplicated cauliflower mosaic virus 35S promoter and a t-nos terminator sequence and O1K consisting of fused VP1 - 4 coding sequences with VP1 fused to a 6xHis coding sequence (P1-polyprotein). (B) Schematic representation of GOR T-DNA used for agroinfiltration. Binary vector pGK2 (GOR) harbouring the GOR gene. RB and LB refers to the right and left border, respectively. *NptII* refers to the neomycin phosphotransferase gene which confers kanamycin resistance.

Results

Protease sensitivity of model recombinant proteins

Since VP1 was used in the study as one of the model recombinant proteins, the VP1 protein was first treated with either a cysteine (papain) or serine (trypsin) protease to determine VP1 sensitivity to protease treatment (Figure 2). Both proteases degraded VP1 when determined by SDS-PAGE analysis, but with more severe VP1 degradation occurring when treated with

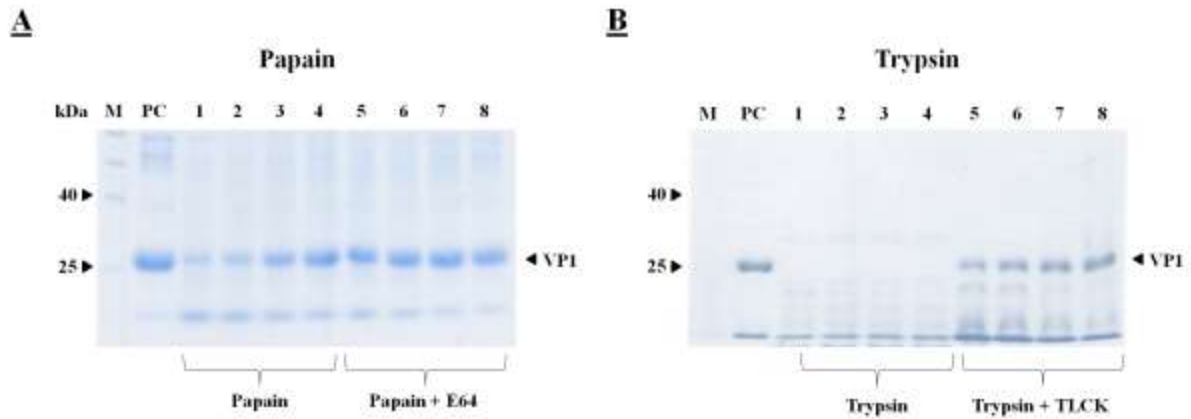


Figure 2 SDS-PAGE analysis of VP1 susceptibility to papain (A) and trypsin (B). Samples were analysed on a 12% SDS-polyacrylamide gel. All sample lanes contain 44 µg of VP1 treated with either papain or trypsin. Lanes 1 to 4 contain either papain (0.2 µg, 0.3 µg, 0.4 µg & 0.5 µg) or trypsin (0.2 µg, 0.3 µg, 0.4 µg and 0.5 µg). Lanes 5 to 8 contain VP1 treated with papain (0.2 µg, 0.3 µg, 0.4 µg and 0.5 µg) together with the cysteine protease inhibitor E64 (1 µM) or VP1 treated trypsin (0.2 µg, 0.3 µg, 0.4 µg and 0.5 µg respectively) together with the serine protease inhibitor TLCK (40 mM). PC (44 µg purified VP1) represents the positive control and M represents a pre-stained protein ladder.

trypsin (Figure 2b). Less degradation occurred when either E64, a cysteine protease inhibitor, or TLCK, a trypsin inhibitor, was added to the reaction mixture (Figure 2a, b). In order to also investigate the influence of proteases *ex planta*, VP1 was further treated using tobacco extracts with and without the addition of a protease inhibitor (Figure 3). VP1 band intensity changed indicating possible protease action; also, some smaller sized bands cross-reacted with the His-antiserum possibly indicating some proteolytic degradation products (Figure 3a). Furthermore, a tobacco extract containing a protease inhibitor resulted in less VP1 degradation (Figure 3b, c).

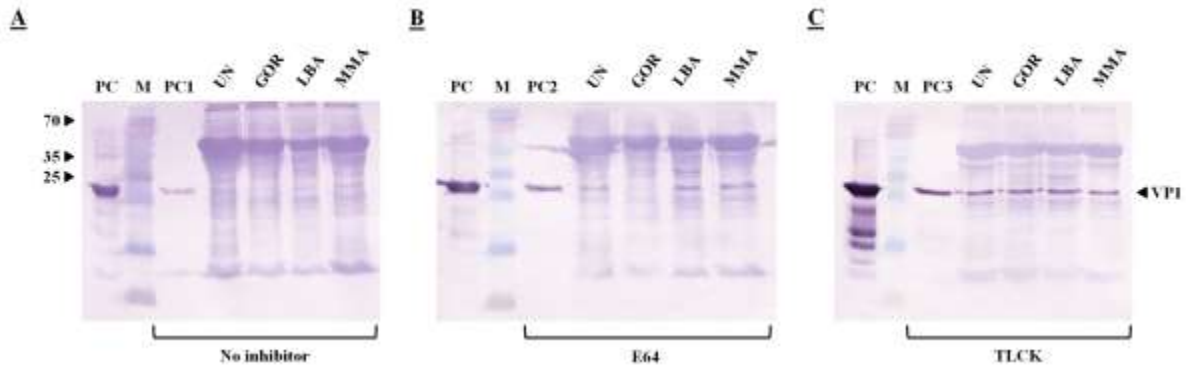


Figure 3 Western blot analysis of VP1 susceptibility with Anti-His antiserum of exogenous VP1 incubated with *N. benthamiana* leaf extracts with or without protease inhibitors. Lanes 1-4 represent 15 µg of VP1 treated with a tobacco extract (230 µg) with either no protease inhibitor (**A**), a cysteine protease inhibitor E64 (100 µM) (**B**) or a serine protease inhibitor TLCK (100 µM) (**C**) added to Arakawa extraction buffer. Lanes UN, GOR, LBA and MMA represent VP1 (15 µg) incubated with tobacco extract only, GOR agroinfiltrated extract, LBA agroinfiltrated extract and MMA agroinfiltrated extract, respectively. Lane M represents a PAGERuler pre-stained protein ladder. PC represents VP1 positive control (60 µg). PC1, PC2 and PC3 represent 15 µg of VP1 incubated with Arakawa extraction buffer.

Protease cleavage sites

In a second step we investigated *in vitro* if model recombinant proteins VP1, GOR as well as P1-polyprotein (VP2 - 4) have proteolytic cleavage sites (Table 1). Subsite nomenclature was assumed from a model created by Schechter and Berger (1967, 1968) where the amino acid residues in a substrate undergoing proteolytic cleavage are designated P1, P2, P3, P4, etc., in the N-terminal direction from the cleaved bond. VP1 was particularly susceptible to papain cleavage with three papain amino acid sites (C²⁵, H¹⁵⁹, N¹⁷⁵) involved in the interaction. When amino acid T (threonine) was used in the P1 substrate position, 62 cleavage sites were found in the polyprotein (Table 1, highlighted in yellow). Within the GOR sequence, 34 cleavage sites were found when amino acid G (glycine) was used in the P1 substrate position (Table 1, highlighted in yellow). Cleavages sites with amino acid W (tryptophan) in the P1

substrate position were, however, not highly abundant in either the P1-polyprotein or GOR sequences with only 3 and 2 sites, respectively, being present (Table 1, highlighted in turquoise). Papain cleaves at TL²¹³ at the end of the VP1 sequence as well as AE²²⁰ at the end of the VP3 sequence whilst cathepsin L cleaves at KE²¹⁸ at the end of the VP2 sequence (Table 3.1, highlighted in grey). Cathepsin H further cleaves VP1 at L²¹³ and at R⁴⁷⁸ of the GOR sequence (Table 1, highlighted in grey). Cathepsin H cleavage sites, with K and L at the P1 substrate position, were further abundant for GOR and VP1 sequences with 34 and 33 sites, respectively (Table 1, highlighted in yellow). VP1 seems to be particularly susceptible to cathepsin H cleavage with a total of 61 potential cleavage sites (Table 1, highlighted in yellow).

To further determine VP1 susceptibility to papain, or trypsin action, ZDOCK protein modelling was then applied between VP1 and the two proteases and the distances (in Ångstroms) between interacting VP1 and protease residues (VP1 cleavage site R⁶⁷ for papain and VP1 cleavage site R²⁶ for trypsin) were determined. The distance in Ångstroms (Å) decreased when VP1 was not modelled together with the additional capsid proteins (VP2, VP3 and VP4) indicating higher VP1 sensitivity to protease action (Figure 4, bottom panels). In contrast, when all other binding sites for VP2, VP3 and VP4 within the P1-polyprotein were permitted in the interaction model with the protease, the distance increased with a weaker VP1-protease interaction and better stability of VP1 (Figure 4, top panels).

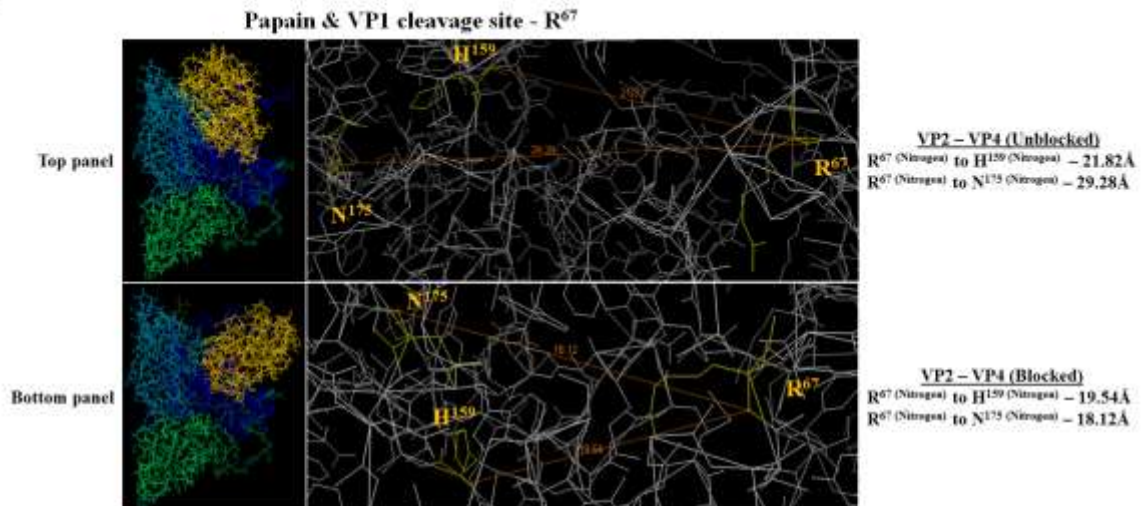
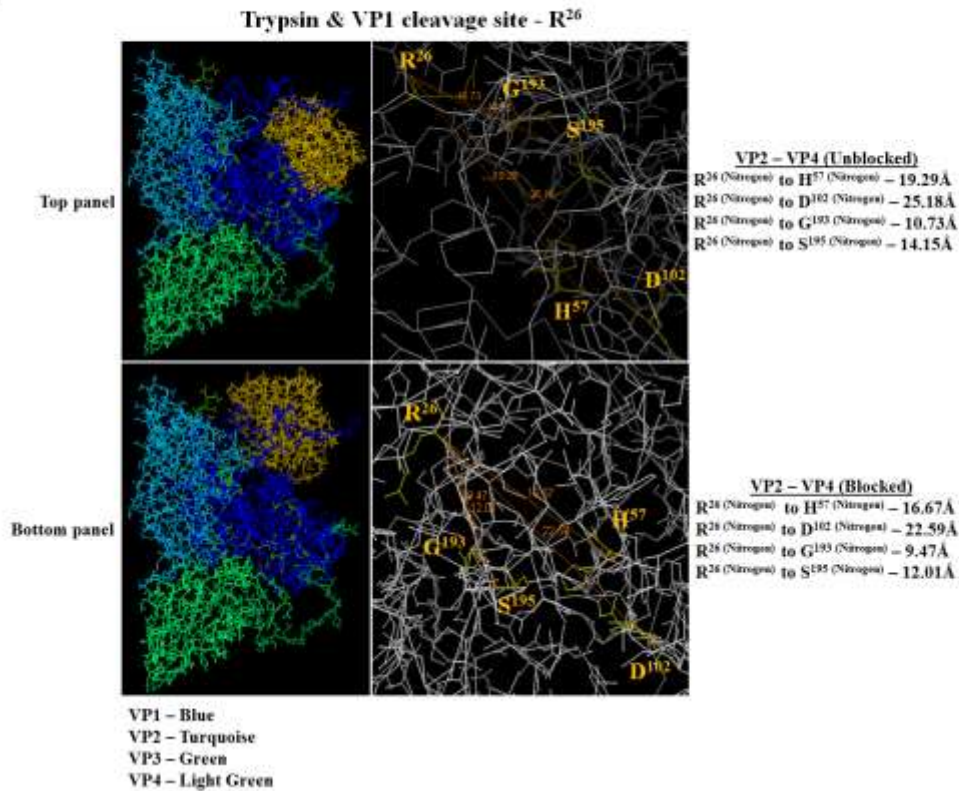
A**B**

Figure 4 Protein docking model of single VP1 with papain and trypsin. **(A)** Docking model of VP1 (cleavage site - R⁶⁷) and papain. **(B)** Docking model of VP1 (cleavage site - R²⁶) and trypsin. Unblocked means all other binding sites (capsid proteins VP2, VP3 and VP4) within the P1-polyprotein were permitted in the interaction model. Blocked means all other binding sites (capsid proteins VP2, VP3 and VP4) within the P1-polyprotein were blocked in the interaction model. Distances between interacting residues are given in Ångstroms (Å).

Transcriptome analysis of agroinfiltrated tobacco leaves

Since we found no significant statistical differences in proteolytic activities between experimental (LBA, P1, GOR) and control (UN & MMA) groups when various cysteine protease activities (cathepsin L and H, legumain proteolytic activities) were measured with fluorogenic substrates 24 hr post infection (pi) (Figure 5), possible expression of proteases due to agroinfiltration was also investigated using RNA-seq analysis. When a limit of at least two-times higher protease expression in agroinfiltrated tissues based on FPKM (fragments per kilobase of transcript per million mapped reads) values was set, the most expressed proteases were cysteine proteases (7 proteases) belonging to 3 different cysteine protease families (C1, C13 and C14) (Figure 6, Table 2) with most cysteine proteases belonging to the C1 family of cysteine proteases (4 proteases). Further, other proteases expressed as a consequence of agroinfiltration belonged to metallo-proteases (2 proteases) in 2 families (M38 and M67), aspartic proteases (2 proteases) in 2 families (A1 and A22) and threonine proteases (1 protease) in 1 family (T1) (Table 2). In comparison to all other proteases, the cysteine protease RD21 ([XM_009614860.1](#)) was the highest expressed cysteine protease. Agroinfiltration increased expression of this cysteine protease about 4-times (Table 2). The protease most induced was *NbVPE-1b*, a vacuolar processing enzyme, belonging to the C13 cysteine protease family (Table 2). This protease was expressed about 95-times higher (based on FPKM-values) in agroinfiltrated leaves than in non-infiltrated leaves. Such a high expression was also found when tobacco leaves were infiltrated with a construct to produce the P1-polyprotein (75-times) and with a construct for GOR production (92-times). For all other identified proteases belonging to different classes, expressions in non-treated leaves were much lower compared to agroinfiltration-induced expression which was in the range of

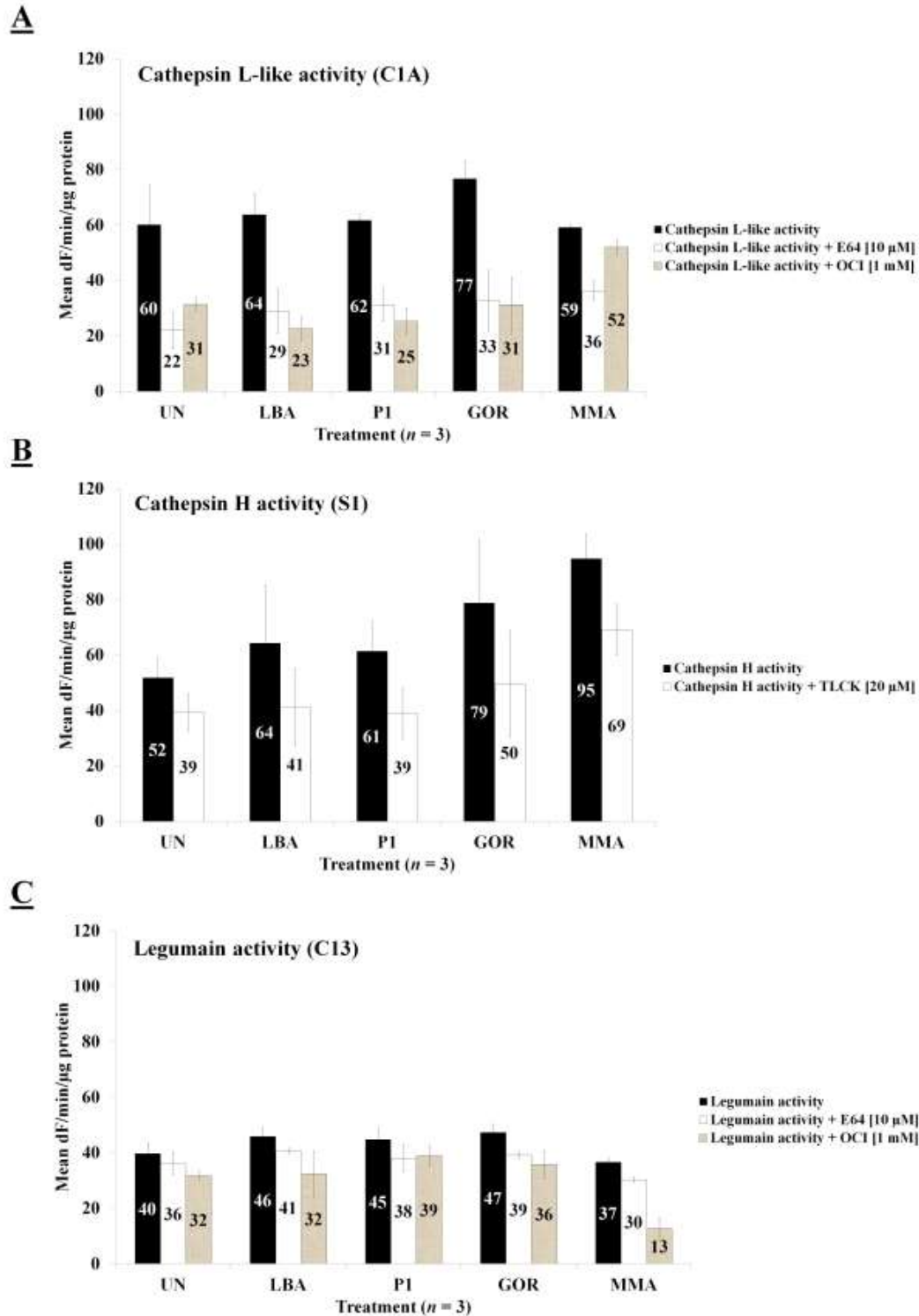


Figure 5 Protease activities of cathepsin L-like, cathepsin H-like and legumain-like in control (UN & MMA) and experimental groups (LBA, P1 and GOR). (A) Cathepsin L-like activity, (B) cathepsin-H like activity, (C) legumain-like activity. The y-axis represents the mean activities expressed as fluorescence units (dF) per min per μg protein. Mean activities of three biological replicates are shown within bars. Error bars indicate standard error of the mean (SEM).

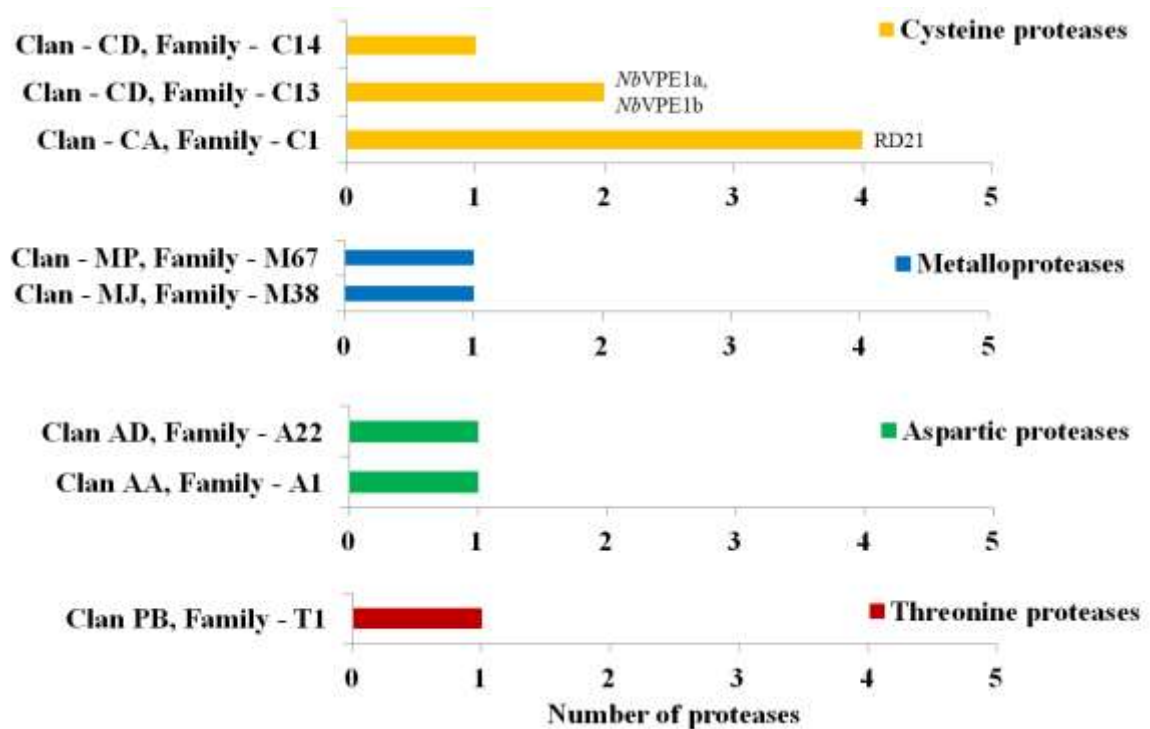


Figure 6 Proteolytic class distribution bearing the clans and families of the five major classes of proteases in *N. benthamiana* from the MEROPS database found in RNA-Seq datasets and number of protease(s) identified in various classes to be expressed at 2-times higher based on FPKM values after agroinfiltration. Genes further validated by RT-qPCR (RD21 and *NbVPE*). Yellow – Cysteine proteases. Blue – Metalloproteases. Green – Aspartic proteases. Maroon – Threonine proteases.

2-30-times more. This increase was irrespective of infiltration with *Agrobacterium* alone or with a construct allowing P1-polyprotein or GOR expression.

Protease reverse transcription quantitative real-time PCR (RT-qPCR)

RT-qPCR was carried out to confirm the RNA-Seq data. Three sequences coding for cysteine protease-like proteins (RD21, *NbVPE-1a* and *NbVPE-1b*) were selected for confirmation. Significantly elevated gene expression was found in all experimental groups (LBA, P1 and GOR) relative to the control (UN), which was set at 1 (Figure 7), in general accordance with the RNA-Seq data as shown by the line graphs of the three respective genes. There were

statistically significant differences ($p < 0.05$) between the expression levels of the different genes assayed for (RD21, *NbVPE-1a*, *NbVPE-1b*). The highest fold increase due to agroinfiltration was also found for, *NbVPE-1b*, comparable to the RNA-Seq result. The discrepancy between the RNAseq and RT-qPCR data for *NbVPE-1b* can be accounted for by the reduced coverage within the RNA-seq datasets compared to the RT-qPCR data.

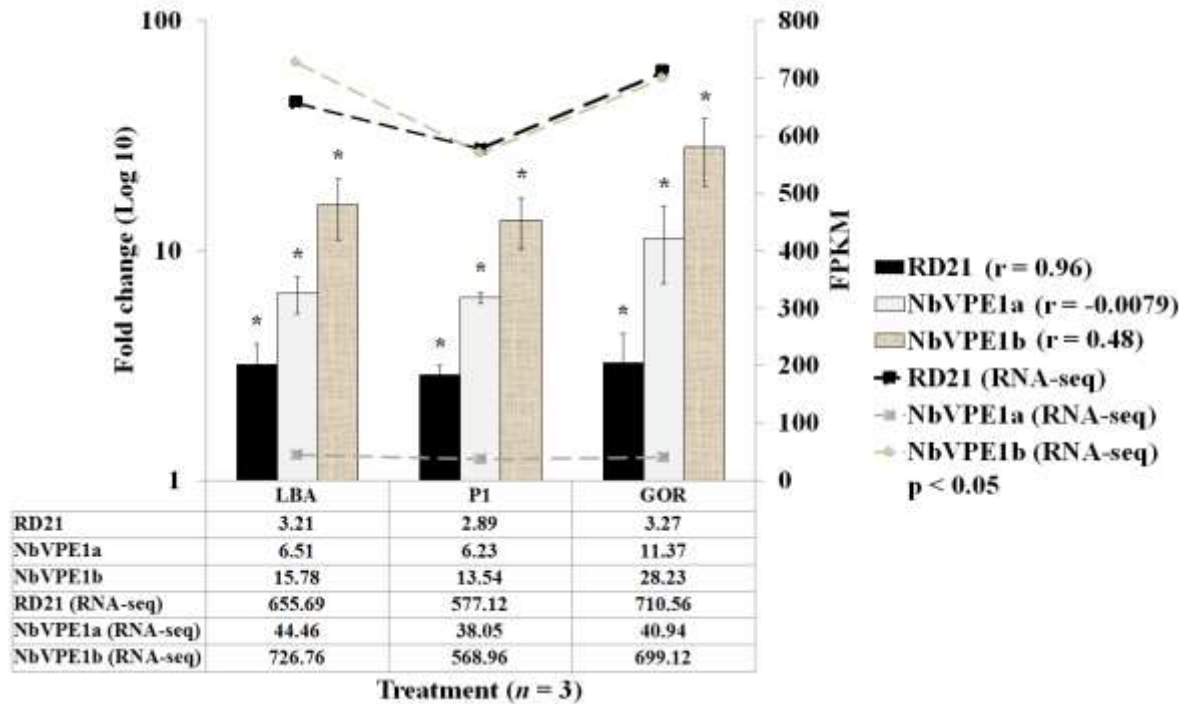


Figure 7 Fold change and FPKM of gene expression of RD21 (black bars), *NbVPE1a* (grey bars) and *NbVPE1b* (canvas bars) in LBA, P1 and GOR experimental groups normalized to reference genes and relative to untreated control (UN) set at 1. Fold changes are represented on left y-axis as log 10 and samples are represented on x-axis. Error bars indicate standard error of mean (SEM) across three biological replicates. FPKM values of RD21 (black line), *NbVPE1a* (grey line) and *NbVPE1b* (canvas line) are represented on the right y-axis. ‘r’ refers to the Pearson correlation coefficient between RNA-seq and RT-qPCR gene expression measurements. Statistically significant differences between gene expression levels (RD21, *NbVPE1a* and *NbVPE1b*) in treatments (LBA, P1 and GOR) compared to control (UN) were determined by two-factor ANOVA with replication (p -value < 0.001) and Bonferroni corrected post-hoc t-tests (p -value < 0.05) are represented as asterisks above graphs (*).

Discussion

The effect of proteases on recombinant protein stability in plant-based expression systems is a field of growing interest. Results in this study have clearly demonstrated that three selected recombinant proteins (P1-polyprotein, VP1 and GOR) are sensitive to protease action. Such action has also been recently associated with lower antigenicity when heterologous proteins were expressed in a plant host.^{13,21} Sensitivity of VP1 against proteases was further confirmed when treated with a protease-containing tobacco extract and when *in silico* proteolytic site analysis was carried out. In particular, cathepsin H-like cysteine proteases might play a major role in VP1 degradation. Our study also provided evidence that proteins, like VP1, are less susceptible to protease degradation when part of a polyprotein.²⁴ Such greater VP1 stability and better processing when expressed in tobacco as a P1-polyprotein, has previously been found.²⁵ Less protease sensitivity was also identified in our protein docking experiments when the distance between interacting residues of VP1 and proteases increased as a result of being part of a polyprotein.

Besides establishing protease sensitivity of model recombinant proteins, sequencing of RNA (RNA-Seq) was applied in our study as a powerful technique to identify any possible proteases expressed during the agroinfiltration process which might compromise recombinant protein production. RNA-Seq was conceived about ten years ago and has become a preferred technology for transcriptome and gene analysis with a number of advantages.^{26,27} RNA-Seq is not dependent on sequence knowledge being available *a priori* and provides a direct measure of RNA abundance. The technology, although powerful, still has the problem of introducing a certain degree of bias due to the sequencing of pooled samples. Rigorous post-sequencing bioinformatics data analysis, as done in our study, is therefore required as well as a final

validation of determined gene expression via alternative methods such as RT-qPCR (reverse transcription quantitative real-time PCR). However, a problem faced in our study, when working with tobacco, was mismatches between gene identifiers and existing gene annotations rendering the annotation process an arduous and challenging task.

By performing RNA-seq analyses, we were able to obtain a transcriptomic profile for proteases in operation during the early stages of agroinfiltration. Members of the C1 and C13 cysteine protease families were more expressed in agroinfiltrated leaves. This result supports our original hypothesis that proteases and, in particular, cysteine proteases are expressed during agroinfiltration. This is consistent with previous findings that C1 proteases are specifically up-regulated in response to pathogenic microbes.^{19,28} Studies have also shown that C1 family RD21-like cysteine proteases, [XM_009614860.1](#) in our study, are important components of plant immunity.^{29,30} These proteases are also induced in senescing leaves in conjunction with the induction of vacuolar processing enzymes (VPEs).³¹ Endogenous cysteine protease inhibitors might further be possible interaction partners of RD21-like proteases to prevent RD21 activity. Our previous finding of better stability of a recombinant protein (GOR) produced via agroinfiltration in tobacco leaves engineered with the rice cysteine proteases inhibitor OC-I supports the idea that RD21-like proteases might have been inhibited by OC-I expression in a transgenic tobacco leaf resulting in better GOR production.²³

A further C1 cysteine protease (ALP; [XM_009594290.1](#)) with aleurain-like activity was found to be expressed 4-times more. Plant aleurain, first isolated from barley and localized in the vacuole, is an amino-peptidase with a number of similarities to animal cathepsin H. These similarities include heterogeneity of charge forms, position of the NH₂-terminus of the

mature protein, and a similar pH-activity profile.³² Several cleavage sites were found in this study for cathepsin H-like proteases in our model recombinant proteins possibly compromising their stability. *N. tabacum* contains cathepsin H-like proteases, such as *NtCP-23* and *NtCP1*, with higher expression during natural senescence.^{23,33,34}

The most prominent cysteine proteases induced in our study by agroinfiltration were two VPEs (*NbVPE-1a* & *NbVPE-1b*; [AB181187.1](#) & [AB181188.1](#)). VPEs, so far, have not been extensively investigated in the context of degrading recombinant proteins. However, it is already known that prolonged incubation of anti-HIV antibodies 2F5 with high amounts of a VPE results in the formation of a 30-kDa degradation product implicating the involvement of VPEs in the degradation of this antibody.²¹ Such degradation might also be relevant for VP1 degradation in our study. VPE cleavage sites were identified in *in silico* cleavage analyses of VP1 and P1-polyprotein as well as that of GOR. Vegetative-type VPEs are generally expressed during senescence or the pathogen-induced hypersensitive response (HR).³⁵ VPEs belong to the CD clan of cysteine proteases and within the clade they form part of the C13 family.³⁶ VPEs contribute to the senescence process and PCD (programmed cell death) by participating in the collapse of the vacuole membrane with the release of proteases into the cell.³⁷ VPEs are asparaginyl endopeptidases cleaving peptide bonds on the C-terminal side of asparagine (Asn) and aspartic (Asp) residues of pro-protein precursors to generate mature proteins.³⁸ They occur along the secretory pathway and are located in the vacuole, except for a single cell wall specific VPE.^{15,39} Although VPEs have caspase-1 activity, they are structurally unrelated to caspases. Plants have evolved their own unique alternative regulated cellular suicide strategy that differs from animals with VPEs located in the vacuole as opposed to the cytosol.⁴⁰

In our study, other possible protease candidates were also identified which might compromise recombinant protein production. This included the C1 cysteine proteases [XM_009798142.1](#), a xylem cysteine peptidase 2 (XCP2) as well [XM_009772991.1](#), a xylem bark cysteine peptidase 3 (XBCP3). Both proteases are involved in various cellular proteolytic processes. Their expression and induction was, however, not comparable to [XM_009614860.1](#) (RD21) and any importance of these two proteases in recombinant protein degradation has yet to be determined. Also the expression of a metacaspase-1 ([XM_009795875.1](#)), involved in apoptosis, was identified in our study, as well as various aspartic, threonine and metallo-proteases. Serine and metallo-proteases have been already found in extracellular root and plant cell culture media. Lallemand *et al.* (2015) recently proposed that they are prime candidates in compromising protein stability. Since serine- or metallo-proteases were not strongly expressed in our study following agroinfiltration, their importance in compromising recombinant protein stability requires more detailed investigation.

Overall, our study has provided a more detailed insight of protease sensitivity of recombinant proteins, in particular, against cysteine proteases. Through, transcriptomic profiling and gene expression analysis, further evidence was provided that several cysteine proteases including proteases with RD21-like and aleurain-like activity as well as VPEs are up-regulated during the early stages of agroinfiltration and might compromise recombinant protein production. However, the identified proteases still require more detailed analyses to confirm their direct involvement in recombinant protein degradation. Ultimately, characterization of identified proteases might also contribute towards establishing a protease library to be screened before any recombinant protein production is envisaged via agroinfiltration.

Methods

Plant material and growth conditions

N. benthamiana seeds were obtained from Dr Ereck Chakauya (CSIR, Pretoria, South Africa) and were germinated in plastic trays in germination mix. Seedlings were grown at a 12/12 hours light/dark cycle with a day/night temperature of 26°C/20°C and 80% (v/v) relative humidity in a growth chamber (Sanyo, Bensenville, USA). Plants were grown for 12 weeks in order to obtain fully expanded leaves suitable for agroinfiltration.

VP1 expression and purification

Recombinant N-terminal His-tagged VP1 was expressed in *E. coli* M15 cells with bacteria grown in LB medium containing 100 µg mL⁻¹ ampicillin, at 37°C up to a cell mass of 0.5 (OD₆₀₀).⁴¹ Expression was induced with 1 mM IPTG for 5 h at 37°C and cells were harvested by centrifugation and stored at -20°C until use. All purification procedures were carried out as previously described with a purification column (Bio-Rad, CA, USA) and elution with imidazole.⁴²

Protein extraction

Whole leaf proteins were extracted in ice-cold Arakawa buffer containing 10 mM L-cysteine, 200 mM Tris-HCl, pH 8.0, 100 mM NaCl, 400 mM Sucrose, 10 mM EDTA, 14 mM 2-Mercaptoethanol, 0.05% Tween-20.⁴³ Total soluble protein (TSP) amount was determined

with a commercial protein determination kit and standardized across samples within experiments (Bio-Rad, Hercules, CA).

VP1 treatment

Varying amounts (0.5, 0.4, 0.3, and 0.2 μg) of either the cysteine protease, papain (Sigma-Aldrich, Germany) or the serine protease, trypsin (Sigma-Aldrich, Germany) were added to 44 μg of purified VP1 protein in 20 μL sodium phosphate buffer (pH 6.0) and incubated for 5 min at 37°C. Tobacco extracts (230 μg total protein) were added to 15 μg of purified VP1 protein and incubated first for 2 hr at 25°C and then for 2 hr at 37°C.

SDS-PAGE and immunoblotting

Protein samples were boiled at 95°C for 5 min in a 4x SDS-containing reducing sample buffer. VP1 stability against proteases was fractionated on a 15% SDS-PAGE gel under reducing conditions with a Mini-PROTEAN® Electrophoresis System (Bio-Rad, Hercules, CA).⁴⁴ Western-blotting was carried out with a 0.45 μm Nitrocellulose membrane (Bio-Rad, Hercules, CA) and a Mini-Trans-Blot electrophoretic transfer cell for protein transfer (Bio-Rad, Hercules, CA) at a constant current of 300 mA for 3 hr. For blotting, membranes were blocked overnight in 5% (w/v) skim milk solution in TBST buffer containing 100 mM Tris, 154 mM NaCl, 0.1% Tween-20, pH 7.5. Incubation with the primary antibody (1:7000 dilution of the His-antiserum raised in rabbit) was done overnight in a 5% (w/v) skim milk-TBST solution. Incubation with the secondary antibody (1:10000 dilution of goat anti-rabbit IgG conjugated with alkaline phosphatase; AbD Serotec, UK) was then done for 1 hr.

Membranes were finally developed with the AP (Alkaline Phosphatase) Conjugate Substrate Kit as described in the manufacturers protocol (Bio-Rad, Hercules, CA).

Protease activity measurement

Total soluble protein (TSP) from foliar extracts (36 µg) were used for measuring cathepsin L-like protease activity in a 50 mM sodium phosphate buffer, pH 6.0 (Sigma-Aldrich, Germany) containing 10 mM L-cysteine (Sigma-Aldrich, Germany) and cathepsin H-like protease activity in a 50 mM Tris buffer (pH 6.0). Extracts were mixed and transferred into black, flat-bottom polysorp 96 well plates (Nunc, AEC Amersham) for measuring fluorescence. Before measuring fluorescence, 8 µM of the substrate Z-Phe-Arg-7-amido-4-methylcoumarin hydrochloride (Z-Phe-Arg-MCA, Sigma-Aldrich, Germany) for cathepsin L-like activity, or Arginine-7-amido-4-methylcoumarin hydrochloride (Arg-NMec HCl, Sigma-Aldrich, Germany) for cathepsin H-like activity, was added in a final volume of 100 µL. Activity was measured kinetically over a 10 min time period with 20 seconds (sec) of shaking before the first cycle. Fluorescence development was measured with a fluorometer (BMG FluoStar Galaxy, Germany) at 25°C with excitation and emission wavelengths of 360 nm and 450 nm, respectively. Legumain activity in TSP (36 µg) was measured in a legumain assay buffer, pH 5.8, containing 1 mM DTT, 39.5 mM citric acid, 121 mM Na₂HPO₄, 1 mM Na₂EDTA, 0.01% CHAPS (Sigma-Aldrich, Germany). Extracts were mixed and prepared as described above and protease activity was measured fluorometrically after addition of 1 mM of substrate Z-Ala-Ala-Asn-AMC (Z-AAN-AMC, Bachem, Germany) as described above for cathepsin L and H activity. A broad spectrum commercial inhibitor N-[N-(L-3-transcarboxyirane-2-carbonyl)-L-Leucyl]-agmatine (E-64, Sigma-Aldrich, Germany) or 10 µM. purified oryzacystatin-I (OCI) (Department of Plant Sciences, University of Pretoria,

South Africa) was used to inhibit cathepsin L-like and legumain protease activity. A commercial inhibitor N-a-tosyl-L-lysine chloromethyl ketone hydrochloride (TLCK, Sigma-Aldrich, Germany) was used to inhibit cathepsin H-like protease activity at 20 μ M.

***In silico* proteolytic cleavage analyses**

Proteolytic cleavage assays were conducted *in silico* with CLC Main Workbench 6.6.1 (<http://www.clcbio.com>) based on various proteases substrate specificities and the Schechter and Berger (1967, 1968) subsite nomenclature.⁴⁵⁻⁴⁹ Protein structures and amino acid sequences were obtained from the RCSB Protein Data Bank (PDB) (www.rcsb.org).⁵⁰ Structures obtained were: P1 – PDB ID: 1FOD⁵¹, GOR – PDB ID: 3GRS⁵², Papain – PDB ID: 9PAP⁵³ and Trypsin – PDB ID: 1UTN.⁵⁴ Substrate specificity profiling was determined using guidelines as previously described.⁴⁷ Protein modelling was conducted between VP1 and papain as well as trypsin, by applying ZDOCK (<http://zdock.umassmed.edu/>).^{55,56} Both the unblocked and blocked settings were used on binding sites within the other P1-polyprotein chains (VP2, VP3 and VP4) when conducting the modelling. Blocking binding sites in the other chains were conducted to simulate an interaction with only VP1 and the respective proteases. Models were visualized in 3D-Mol Viewer (a component of Vector NTI 9.1.0, Invitrogen) and distances between interacting residues were measured using the measure distance tool. The setting structure was used as the colour theme. VP2 – 4 were hidden from selections to highlight the interaction between VP1 and the protease. Distances were measured in the measure distance mode in Ångstroms (Å) between defined molecules in amino acids.

Agroinfiltration

Syringe agroinfiltration was used to infiltrate the tobacco leaf surface.⁷ Cultures were maintained in lysogeny broth (LB) medium supplemented with 50 $\mu\text{g mL}^{-1}$ kanamycin and 50 $\mu\text{g mL}^{-1}$ rifampicin. For agroinfiltration, bacteria were grown to the stable phase at 28°C to an OD_{600} of 1 and collected by centrifugation at 4000g. Bacterial pellets were re-suspended in MMA medium (10 mM 2-[N-morpholino] ethanesulfonic acid] (MES) buffer, pH 5.6, containing 100 μM acetosyringone and 10 mM MgCl_2). For RNA-seq (RNA sequencing) analyses, the first fully-expanded leaves of the upper four individual leaves were infiltrated as previously described using a needle-less syringe with the vector pB+O1KP1-*HT* (gift from Prof George Lommonossoff; John Innes Centre, Norwich, UK) containing the P1-polyprotein (P1) coding sequence (Figure 1a) or with the pGK2 construct (Figure 1b) containing the glutathione reductase (GOR) coding sequence.^{7,19} Leaves were also agroinfiltrated with *A. tumefaciens* strain LBA4404 (LBA) alone. Uninfiltrated leaf material (UN) and infiltration with MMA medium were applied as negative controls to avoid confounding effects due to experimental conditions. Infiltrated plants were kept in a growth cabinet (Sanyo, Bensenville, USA) for 24 hr. Three biological (plant) replicates were used for each treatment allowing for statistical treatment of data. Plants were kept in an environmentally controlled growth room and watered daily. Four leaves per plant were harvested after 24 hr, as source material for subsequent protein and RNA extraction. Leaf samples were frozen immediately in liquid nitrogen and stored at -80°C until protein or RNA extraction was carried out.

RNA extraction

RNA extraction on leaf tissue samples was carried out with the Trizol method by employing macro-dissection.⁵⁷ Ribolock (Thermo Fisher Scientific, Waltham, MA USA) was added to the final volume of RNA in a ratio of 1:10 (v/v). On-column DNase (Thermo Fisher Scientific, Waltham, MA USA) digests were performed with the Qiagen RNeasy Mini Kit (Qiagen, Valencia, CA, USA). RNA samples were initially analysed on a full-spectrum spectrophotometer Nanodrop® ND-1000 (Thermo Fisher Scientific, Waltham, MA USA) and subsequently sent for RNA quality analyses on the Experion™ Automated Electrophoresis System (Bio-Rad, CA, USA). After quantification, RNA samples from each sample were equivalently pooled for RNA-seq analyses. Biological replicates for each group were kept separate for RT-qPCR (reverse transcription quantitative real-time PCR) analyses.

RNA-seq and mapping

A transcriptomic library was constructed from paired end reads of ~90 bp in size which were generated from the HiSeq 2000 sequencing system (Illumina® sequencing, San Diego, USA) at the Beijing Genomics Institute (BGI Tech Solutions Co., Ltd, Hong Kong, China). Galaxy (Department of Bioinformatics, University of Pretoria), a platform for working with sequencing data, was applied to visualize, interpret, and conduct further analyses on the data generated.⁵⁸⁻⁶⁰ As part of the filtering process, reads with adaptors, unknown nucleotides larger than 5% and with low quality (more than 20% of the bases' qualities are less than 10 in a read) were removed (BGI Tech Solutions Co., Ltd, Hong Kong, China). The FastQC tool was applied to perform QC checks on the data (Supplementary Figure 1). The FastQ groomer was then used to convert the data into a format amenable for subsequent interpretation.⁶¹ The

draft assembly of the *N. benthamiana* genome was applied as a reference genome.^{2,3,62} With TopHat2 RNA-Seq reads were aligned to the tobacco genome available from the Solanaceae Genomics Network (SGN) at url (ftp://ftp.solgenomics.net/genomes/Nicotiana_benthamiana/).^{63,64} Tophat2 was carried out with a mean inner distance between mate pairs of 120 and a standard deviation of 30.

Transcript abundances and gene annotation

Cufflinks was applied performing bias correction and in default mode for all other parameters to assemble transcripts and estimate abundances.⁶⁵ High-confidence transcripts were obtained from identified transcripts (i.e., transcripts with FPKM value in the case of cufflinks > 0) by filtering for a FPKM 95 % confidence interval lower boundary greater than zero and FPKM value \geq 0.001. With MEROPS, Sol Genomics Network (http://solgenomics.net/organism/Nicotiana_benthamiana/genome), and TAIR databases, proteases were mined from datasets.^{2,3,62,66,67} Blast analyses against the draft genome of *N. benthamiana* were conducted with the BLASTN tool.^{2,3,62} Tracking IDs were obtained for transcripts of interest and applying the Log10 of FPKM expression data values and transcript abundances were established.⁶⁸ Gene annotation was conducted using Blast2GO, TAIR as well as NCBI databases.^{67,69,70}

RT-qPCR analyses

mRNA transcripts for three proteolytic candidate proteases (RD21, *NbVPE1a* and *NbVPE1b*) were assayed by RT-qPCR with a Bio-Rad CFX C1000™ Real-Time PCR Detection System (Bio-Rad, CA, USA). cDNA synthesis was done with the Promega GoScript™ Reverse

Transcription System (Madison, Wisconsin, USA). RT-qPCR assays were carried out in accordance with the MIQE guidelines and optimized for annealing temperature and primer efficiency.⁷¹ Reactions contained 10 μ M forward and reverse primer and 2.5 μ L of cDNA template. Supplementary Table 1 provides information of primer sets used. No-template mixture controls were included in each 96-well plate. Thermocycling parameters included initial denaturation at 95°C for 3 min; 39 cycles of denaturation at 95°C for 10 sec, annealing and extension at the abovementioned ranges for 30 sec, denaturation at 95°C for 10 sec. Melt curve analyses were performed thereafter from 65°C to 95°C with 0.5 increments. Fold changes in gene expression were determined with the Livak method.⁷² For comparative purposes, relative gene expression was defined with the value of 1 in control plants.

Statistical analysis

Statistically significant changes in gene expression between control and experimental groups were determined using ANOVA: single/two-factor with replication (p-value < 0.05) applying Microsoft Excel software 2010 version 14 (Microsoft Corporation). Bonferroni corrected post-hoc t-tests (p-value < 0.05) were subsequently performed.

Author contribution statement

PP and KJK conceptualised and designed the experiment. PP conducted the experiments. BJV financially supported the project and provided analytical tools and scientific intellectual input in data interpretation. MEM helped in running enzymatic assays. CAC and SGVW helped in analysing the RNAseq data and manuscript reading. All authors read and approved the manuscript.

Acknowledgments

We thank Professor Yong Suk Jang, Professor Moon Sik Yang and Dr Tae Geum Kim for providing us with the VP1 gene and Dr Huy for assisting with the VP1 purification, Professor George Lomonosoff for providing the P1 construct, and Professor Christine Foyer for providing the GOR construct. We also thank Dr Francois Maree (Onderstepoort Veterinary campus, University of Pretoria, South Africa) who kindly provided the Anti-FMDV polyclonal antiserum and Professor Dominique Michaud for assisting us to establish the agroinfiltration technique. Our research was supported by the National Research Foundation (NRF) and the Genomics Research Institute (GRI), South Africa as well as NRF incentive funding to Karl Kunert and a NRF bursary to Priyen Pillay.

Conflict of interest

The authors declare that they have no conflict of interest.

References

1. Zhang Y, Li D, Jin X, Huang Z. Fighting Ebola with ZMapp: spotlight on plant-made antibody. *Sci China Life Sci* 2014; 57:987-8.
2. Bombarely A, Rosli HG, Vrebalov J, Moffett P, Mueller LA, Martin GB. A draft genome sequence of *Nicotiana benthamiana* to enhance molecular plant-microbe biology research. *Mol Plant Microbe In* 2012; 25:1523-30.
3. Goodin MM, Zaitlin D, Naidu RA, Lommel SA. *Nicotiana benthamiana*: Its history and future as a model for plant-pathogen interactions. *Mol Plant Microbe In* 2008; 21:1015-26.

4. Faino L, de Jonge R, Thomma BP. The transcriptome of *Verticillium dahliae*-infected *Nicotiana benthamiana* determined by deep RNA sequencing. *Plant Signal Behav* 2012; 7:1065.
5. Abiri R, Valdiani A, Maziah M, Shaharuddin NA, Sahebi M, Yusof ZNB, et al. A Critical Review of the Concept of Transgenic Plants: Insights into Pharmaceutical Biotechnology and Molecular Farming. *Curr Issues Mol Biol* 2015; 18:21-42.
6. Wagner B, Fuchs H, Adhami F, Ma Y, Scheiner O, Breiteneder H. Plant virus expression systems for transient production of recombinant allergens in *Nicotiana benthamiana*. *Methods* 2004; 32:227-34.
7. D'Aoust M-A, Lavoie P-O, Belles-Isles J, Bechtold N, Martel M, Vézina L-P. Transient expression of antibodies in plants using syringe agroinfiltration. *Recombinant Proteins From Plants*, 2009:41-50.
8. Hoorn RALVd, Laurent F, Roth R, Wit PJGMD. Agroinfiltration is a versatile tool that facilitates comparative analyses of Avr9/Cf-9-Induced and Avr4/Cf-4-induced necrosis. *Phytopathology* 2000; 13:439 - 46.
9. Outchkourov NS, Rogelj B, Strukelj B, Jongsma MA. Expression of sea anemone equistatin in potato. Effects of plant proteases on heterologous protein production. *Plant Physiol* 2003; 133:379-90.
10. Donini M, Lombardi R, Lonoce C, Di Carli M, Marusic C, Morea V, et al. Antibody proteolysis: a common picture emerging from plants. *Bioengineered* 2015:doi: 10.1080/21655979.2015.1067740.
11. Miletic S, Simpson DJ, Szymanski CM, Deyholos MK, Menassa R. A plant-produced bacteriophage tailspike protein for the control of *Salmonella*. *Front Plant Sci* 2015; 6:1-9.

12. Faye L, Boulaflous A, Benchabane M, Gomord V, Michaud D. Protein modifications in the plant secretory pathway: Current status and practical implications in molecular pharming. *Vaccine* 2005; 23:1770-8.
13. Castilho A, Windwarder M, Gattinger P, Mach L, Strasser R, Altmann F, et al. Proteolytic and N-Glycan Processing of Human α 1-Antitrypsin Expressed in *Nicotiana benthamiana*. *Plant Physiol* 2014; 166:1839-51.
14. Delannoy M, Alves G, Vertommen D, Ma J, Boutry M, Navarre C. Identification of peptidases in *Nicotiana tabacum* leaf intercellular fluid. *Proteomics* 2008; 8:2285-98.
15. Goulet C, Khalf M, Sainsbury F, D'Aoust MA, Michaud D. A protease activity–depleted environment for heterologous proteins migrating towards the leaf cell apoplast. *Plant Biotechnol J* 2012; 10:83-94.
16. Benchabane M, Goulet C, Rivard D, Faye L, Gomord V, Michaud D. Preventing unintended proteolysis in plant protein biofactories. *Plant Biotechnol J* 2008; 6:633-48.
17. Roberts IN, Caputo C, Criado MV, Funk C. Senescence-associated proteases in plants. *Physiol Plant* 2012; 145:130-39.
18. Goulet C, Goulet C, Goulet M-C, Michaud D. 2-DE proteome maps for the leaf apoplast of *Nicotiana benthamiana*. *Proteomics* 2010; 10:2536-44.
19. Robert S, Khalf M, Goulet M-C, D'Aoust M-A, Sainsbury F, Michaud D. Protection of recombinant mammalian antibodies from development-dependent proteolysis in leaves of *Nicotiana benthamiana*. *PloS one* 2013; 8:1-9.
20. Lallemand J, Bouché F, Desiron C, Stautemas J, De Lemos Esteves F, Périlleux C, et al. Extracellular peptidase hunting for improvement of protein production in plant cells and roots. *Front Plant Sci* 2015; 6:1-10.
21. Niemer M, Mehofer U, Torres Acosta JA, Verdianz M, Henkel T, Loos A, et al. The human anti-HIV antibodies 2F5, 2G12, and PG9 differ in their susceptibility to proteolytic

degradation: Down-regulation of endogenous serine and cysteine proteinase activities could improve antibody production in plant-based expression platforms. *Biotechnol J* 2014; 9:493-500.

22. Veena, Jiang H, Doerge R, Gelvin SB. Transfer of T-DNA and Vir proteins to plant cells by *Agrobacterium tumefaciens* induces expression of host genes involved in mediating transformation and suppresses host defense gene expression. *Plant J* 2003; 35:219-36.

23. Pillay P, Kibido T, Plessis M, Vyver C, Beyene G, Vorster BJ, et al. Use of Transgenic Oryzacystatin-I-Expressing Plants Enhances Recombinant Protein Production. *Appl Biochem Biotechnol* 2012; 168:1608-20.

24. Butt TR, Edavettal SC, Hall JP, Mattern MR. SUMO fusion technology for difficult-to-express proteins. *Protein Express Purif* 2005; 43:1-9.

25. Pan L, Zhang Y, Wang Y, Wang B, Wang W, Fang Y, et al. Foliar extracts from transgenic tomato plants expressing the structural polyprotein, P1-2A, and protease, 3C, from foot-and-mouth disease virus elicit a protective response in guinea pigs. *Vet Immunol Immunop* 2008; 121:83-90.

26. Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nat methods* 2008; 5:621-8.

27. Wang Z, Gerstein M, Snyder M. RNA-Seq: a revolutionary tool for transcriptomics. *Nat Rev Genet* 2009; 10:57-63.

28. Gilroy EM, Hein I, Van Der Hoorn R, Boevink PC, Venter E, McLellan H, et al. Involvement of cathepsin B in the plant disease resistance hypersensitive response. *Plant J* 2007; 52:1-13.

29. Gu C, Shabab M, Strasser R, Wolters PJ, Shindo T, Niemer M, et al. Post-translational regulation and trafficking of the granulin-containing protease RD21 of *Arabidopsis thaliana*. *PloS one* 2012; 7:1-11.

30. Shindo T, Misas-Villamil JC, Hörger AC, Song J, van der Hoorn RA. A role in immunity for Arabidopsis cysteine protease RD21, the ortholog of the tomato immune protease C14. *PLoS one* 2012; 7:1-9.
31. Kinoshita T, Yamada K, Hiraiwa N, Kondo M, Nishimura M, Hara-Nishimura I. Vacuolar processing enzyme is up-regulated in the lytic vacuoles of vegetative tissues during senescence and under various stressed conditions. *Plant J* 1999; 19:43-53.
32. Holwerda BC, Rogers JC. Purification and Characterization of Aleurain A Plant Thiol Protease Functionally Homologous to Mammalian Cathepsin H. *Plant Physiol* 1992; 99:848-55.
33. Ueda T, Seo S, Ohashi Y, Hashimoto J. Circadian and senescence-enhanced expression of a tobacco cysteine protease gene. *Plant Mol Biol* 2000; 44:649-57.
34. Beyene G, Foyer CH, Kunert KJ. Two new cysteine proteinases with specific expression patterns in mature and senescent tobacco (*Nicotiana tabacum* L.) leaves. *J Exp Bot* 2006; 57:1431-43.
35. Yamada K, Shimada T, Nishimura M, Hara-Nishimura I. A VPE family supporting various vacuolar functions in plants. *Physiol Plant* 2005; 123:369-75.
36. Mosolov V, Valueva T. Participation of proteolytic enzymes in the interaction of plants with phytopathogenic microorganisms. *Biochemistry (Moscow)* 2006; 71:838-45.
37. Hara-Nishimura I, Hatsugai N, Nakaune S, Kuroyanagi M, Nishimura M. Vacuolar processing enzyme: an executor of plant cell death. *Curr Opin Plant Biol* 2005; 8:404-8.
38. Hara-Nishimura I, Hatsugai N, Nakaune S, Kuroyanagi M, Nishimura M. Vacuolar processing enzyme: an executor of plant cell death. *Curr Opin Plant Biol* 2005; 8:404-8.
39. Müntz K, Blattner FR, Shutov AD. Legumains-a family of asparagine-specific cysteine endopeptidases involved in propeptide processing and protein breakdown in plants. *J Plant Physiol* 2002; 159:1281-93.

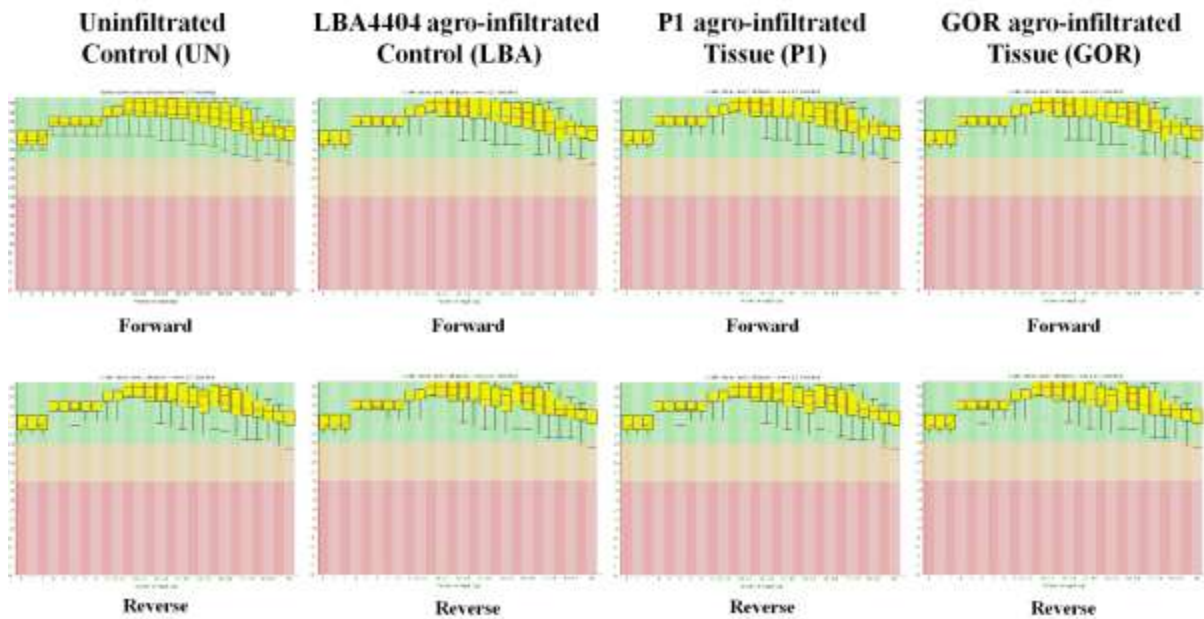
40. Hatsugai N, Kuroyanagi M, Nishimura M, Hara-Nishimura I. A cellular suicide strategy of plants: vacuole-mediated cell death. *Apoptosis* 2006; 11:905-11.
41. Pillay P. Expression of the VP1 antigen from foot-and-mouth disease virus in a bacterial and plant-based expression system. *Plant Sciences*. South Africa: University of Pretoria, 2012:149.
42. Crowe JH, K. *The QIA expressionist*. Chatsworth CA, 1992.
43. Arakawa T, Chong DKX, Merritt JL, Langridge WHR. Expression of cholera toxin B subunit oligomers in transgenic potato plants. *Transgenic Res* 1997; 6:403 - 13.
44. Laemmli UK. Cleavage of structural proteins during the assembly of the head of bacteriophage T4. *Nature* 1970; 227:680 - 5.
45. Schechter I, Berger A. On the size of the active site in proteases. I. Papain. *Biochem Biophys Res Commun* 1967; 27:157-62.
46. Schechter I, Berger A. On the active site of proteases. III. Mapping the active site of papain; specific peptide inhibitors of papain. *Biochem Biophys Res Commun* 1968; 32:898-902.
47. Choe Y, Leonetti F, Greenbaum DC, Lecaille F, Bogyo M, Brömme D, et al. Substrate profiling of cysteine proteases using a combinatorial peptide library identifies functionally unique specificities. *J Biol Chem* 2006; 281:12824-32.
48. Guncar G, Podobnik M, Pungercar J, Strukelj B, Turk V, Turk D. Crystal structure of porcine cathepsin H determined at 2.1 Å resolution: location of the mini-chain C-terminal carboxyl group defines cathepsin H aminopeptidase function. *Structure* 1998; 6.
49. Mathieu MA, Bogyo M, Caffrey CR, Choe Y, Lee J, Chapman H, et al. Substrate specificity of schistosome versus human legumain determined by P1–P3 peptide libraries. *Mol Biochem Parasitol* 2002; 121:99-105.
50. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat T, Weissig H, et al. The protein data bank. *Nucleic Acids Res* 2000; 28:235-42.

51. Logan D, Abu-Ghazaleh R, Blakemore W, Curry S, Jackson T, King A, et al. Structure of a major immunogenic site on foot-and-mouth disease virus. *Nature* 1993; 362:566 - 8.
52. Karplus PA, Schulz GE. Refined structure of glutathione reductase at 1.54 Å resolution. *J Mol Biol* 1987; 195:701-29.
53. Kamphuis IG, Kalk K, Swarte M, Drenth J. Structure of papain refined at 1.65 Å resolution. *J Mol Biol* 1984; 179:233-56.
54. Leiros HKS, Brandsdal BO, Andersen OA, Os V, Leiros I, Helland R, et al. Trypsin specificity as elucidated by LIE calculations, X-ray structures, and association constant measurements. *Protein Sci* 2004; 13:1056-70.
55. Pierce BG, Wiehe K, Hwang H, Kim B-H, Vreven T, Weng Z. ZDOCK server: interactive docking prediction of protein–protein complexes and symmetric multimers. *Bioinformatics* 2014; 30:1771-3.
56. Chen R, Li L, Weng Z. ZDOCK: An initial-stage protein-docking algorithm. *Proteins: Struct Funct Bioinf* 2003; 52:80-7.
57. MacRae E. Extraction of plant RNA. *Protocols for Nucleic Acid Analysis by Nonradioactive Probes*: Springer, 2007:15-24.
58. Goecks J, Nekrutenko A, Taylor J. Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol* 2010; 11:1-13.
59. Giardine B, Riemer C, Hardison RC, Burhans R, Elnitski L, Shah P, et al. Galaxy: a platform for interactive large-scale genome analysis. *Genome Res* 2005; 15:1451-5.
60. Blankenberg D, Kuster GV, Coraor N, Ananda G, Lazarus R, Mangan M, et al. Galaxy: a web-based genome analysis tool for experimentalists. *Curr Protoc Mol Biol* 2010:19.0. 1-.0. 21.

61. Blankenberg D, Gordon A, Von Kuster G, Coraor N, Taylor J, Nekrutenko A. Manipulation of FASTQ data with Galaxy. *Bioinformatics* 2010; 26:1783-5.
62. Knapp S, Chase MW, Clarkson JJ. Nomenclatural changes and a new sectional classification in *Nicotiana* (Solanaceae). *Taxon* 2004; 53:73-82.
63. Bombarely A, Menda N, Teclé IY, Buels RM, Strickler S, Fischer-York T, et al. The Sol Genomics Network (solgenomics.net): growing tomatoes using Perl. *Nucleic Acids Res* 2011; 39:D1149-D55.
64. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol* 2013; 14:1-13.
65. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nature Biotechnol* 2010; 28:511-5.
66. Rawlings ND, Waller M, Barrett AJ, Bateman A. MEROPS: the database of proteolytic enzymes, their substrates and inhibitors. *Nucleic Acids Res* 2013; 42:D503-D9.
67. Lamesch P, Berardini TZ, Li D, Swarbreck D, Wilks C, Sasidharan R, et al. The Arabidopsis Information Resource (TAIR): improved gene annotation and new tools. *Nucleic Acids Res* 2012; 40:D1202-D10.
68. Eisen MB, Spellman PT, Brown PO, Botstein D. Cluster analysis and display of genome-wide expression patterns. *Proc Natl Acad Sci USA* 1998; 95:14863-8.
69. Conesa A, Götz S, García-Gómez JM, Terol J, Talón M, Robles M. Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics research. *Bioinformatics* 2005; 21:3674-6.

70. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 1997; 25:3389-402.
71. Bustin SA, Benes V, Garson JA, Hellemans J, Huggett J, Kubista M, et al. The MIQE guidelines: minimum information for publication of quantitative real-time PCR experiments. *Clin Chem* 2009; 55:611-22.
72. Livak KJ, Schmittgen TD. Analysis of relative gene expression data using real-time quantitative PCR and the $2^{-\Delta\Delta CT}$ method. *Methods* 2001; 25:402-8.

Supplementary Figures



Supplementary Figure 1 Box whisker type plot of quality values across bases at each position in the FastQ files for control (UN) and experimental (LBA, P1 & GOR) sequence data. The elements of the plot are as follows: the central red line is the median value; the yellow box represents the inter-quartile range (25-75%); the upper and lower whiskers represent the 10% and 90% points; the blue line represents the mean quality. The y-axis on the graph indicates the quality scores.