# A Bioinformatics Analysis of Plant Caleosins

Fatemeh Saadat[a*], Houshang Alizadeh[b], Seyed Hadi Razavi[c]

[a]Biotechnology Department, College of Agriculture and Natural Resources,
University of Tehran, Iran
[b]Department of Agronomy and Plant Breeding, University of Tehran, Iran
[c]Department of Food Science and Technology, College of Agriculture and
Natural Resources, University of Tehran, Iran

E-mail:  fa.saadat@chmail.ir
E-mail:  halizade@ut.ac.ir
E-mail:  srazavi@ut.ac.ir

ABSTRACT. Introducing oil bodies to the industry provided an impetus
to know more about the constructive factors. Accordingly, the current
article focused on the description of the plant seed caleosins involved in
oil body formation. For this purpose, bioinformatics analysis was con-
ducted to identify the putative caleosins. Then, the phylogenetic tree
and conserved motifs were extracted from the sequences. Finally, a model
sequence was created by WebLogo to indicate all features including sec-
ondary structure, post-translational modifications and physicochemical
properties. According to the results, the N-terminal region of caleosins
owned more length and negative charges. Although the net charge and
length of domains have changed over evolutionary time, the members
of Brassicaceae had been highly conserved even in the motif sequences.
Considering the importance of these two factors in the construction of
the oil body, natural selection has probably shaped more stable droplets
for lipid storage. However, its confirmation requires more investigation.

**Keywords:** Conserved sequence, Membrane proteins, Nanoparticles, Phy-
logeny.

*Corresponding Author

F. Saadat, H. Alizadeh, S. H. Razavi

**2000 Mathematics subject classification:** 26A33, 45K05, 65T60.

## 1. Introduction

Almost all organisms contain lipid particles made of three main components: fatty acids, phospholipids, and proteins [16]. Despite the lowest portion, proteins are the most effective factors in the formation and function of lipid particles. The plant oil bodies possess three types of proteins: structural proteins (e.g., oleosin and caleosin), enzymes (e.g., lipase), and minimal proteins (e.g., aquaporin) [24]. The oleo-proteins are essential for storage and usage of lipids during the seed germination [27]. Besides, the structural proteins play a key role in constructing artificial oil bodies (AOBs) used in the immobilization, purification, encapsulation, targeting, and bioavailability of recombinant proteins [2].

According to the hydropathy analysis, the structural proteins include a central non-cytoplasmic (N) domain flanked by two cytoplasmic (C) terminal domains (hereafter this structure briefly called CNC). A proline knot motif in the central domain is responsible for mooring into the oil body, and thus supports the terminal domains to settle on the surface in contact with the phospholipids and cell components of the cytoplasm. Although more investigations have been conducted on oleosins, caleosins are the better option to construct the nanoparticles [5; 6] due to the long terminal domains (about 85% of the total sequence) [12]. Moreover, the higher stability of AOBs expected due to the stronger electronegative repulsion of caleosins [4]. The extensive applications of nano-AOBs make it necessary to increase knowledge about the structural proteins to form more stable and efficient oil bodies.

Considering the importance of the structural domains (central and terminal domains) in the construction of oil bodies, a domain analysis was performed to identify the CNC caleosins. Then, the phylogenetic tree and conserved motifs were extracted from the putative caleosins. Finally, a logo was created to describe the physicochemical properties of the protein. According to our knowledge, this is the first study on all CNC caleosins represented in databases.

## 2. Material and Methods

2.1. **Characteristics of plant caleosins.** Protein sequences were collected using caleosin as a keyword in UniProt (http://uniprot.org/). The results were limited to viridiplantae in the taxonomy section. Next, a domain exploration was conducted on the caleosins in Phobius (http://phobius.sbc.su.se) [13]. In this way, the CNC caleosins were selected for further investigation. Then, the charge distribution of each domain was estimated using SAPS (http://www.ebi. ac.uk/Tools/seqstats/saps) [3]. Finally, the total charge, molecular weight (Mw), isoelectric points (pI), and amino acid content were calculated by EMBOSS (http://www.bi.up.ac.za/ EMBOSS).

2.2. **Protein of alignments and phylogenetic studies.** Multiple alignments were performed by ClustalW (http://www.genome.jp/tools/clustalw/) with default parameters. Then, the phylogenetic analysis was carried out using the neighbor-joining (NJ) method implemented in MEGA5.0. The reliability of the tree was assessed by 1000-bootstrap resampling. The maximum likelihood (ML) and parsimony trees were also constructed to confirm the tree topologies [25].

2.3. **Motif analysis of the caleosins.** The caleosin sequences were analyzed by MEME Suite (meme-suite.org/tools/meme) [1] with one occurrence per sequence for five different motifs. The annotations of the motifs were investigated by CDART (http://www.ncbi.nlm.nih.gov/Structure/lexington /lexington.cgi) [10] and GenomeNet (http://www.genome.jp/tools/motif/).

2.4. **Constructing a model for caleosins.** A sequence logo was created by WebLogo (http://weblogo.berkeley.edu/logo.cgi) after removing divergent sequences [8]. Following arranging a sequence model, the secondary structure, post-translational modifications, and ubiquitination were predicted by NPS@: Network Protein Sequence Analysis (https://npsa-prabi.ibcp.fr) [7], CBS Prediction Servers (www.cbs.dtu.dk/services) and UbPred (www.ubpred.org/cgi-bin/ubpred/ubpred.cgi) [19], respectively. Moreover, the hydropathic plot was produced by ProtScale (http://web.expasy.org/protscale/) with Kyte and Doolittle method [14].

## 3. Results and Discussion

3.1. **Characteristics of caleosins.** 100 caleosins with CNC structure were extracted from UniProt data by domain study. The physicochemical properties of the putative proteins are summarized in Table 1.

Further investigation on each domain of CNCs demonstrated a predominantly negative charge on the N-terminal region. However, some caleosins such as K3YU83 from *Setaria italica* mainly carry positive charges. Besides, the N- and C-terminal domains with an average of 100 and 48 residues were the longest and shortest parts of caleosins, respectively. Notably, monocots appear to be more varied in the length of domains. For example, two different caleosins has been detected in *Zea mays* with 30 (A0A1D6JA94) and 165 (B4FKP4) amino acid residues at the N-termini. Nevertheless, there is a balance between the length of N- and C-terminal domains. Thus, a decrease in the N-terminal portion is often compensated by an increase in the C-terminal length.

3.2. **Phylogenetic relationships and evolution of CNCs.** A phylogenetic tree was constructed using the neighbor-joining method (Figure 1). The topology of the tree was highly similar to that of the maximum likelihood and parsimony tree, supporting high confidence in the phylogenetic relationships.

F. Saadat, H. Alizadeh, S. H. Razavi

| | Length | MW(D) | Charge | pI | Tiny$^a$ | Small$^a$ | Aliphatic$^a$ | Aromatic$^a$ | Non-polar$^a$ | Polar$^a$ | Basic$^a$ | Acidic$^a$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Min | 179 | 20849.9 | −6.5 | 4.9 | 20.6 | 42.4 | 14.5 | 9.8 | 26.2 | 40.4 | 11.1 | 8.5 |
| Max | 364 | 41438.7 | 15 | 10.4 | 38.5 | 58.3 | 22.8 | 19.5 | 59.6 | 49.2 | 19.5 | 14.3 |
| Ave. | 255.6 | 28887.4 | 0.9 | 6.8 | 27.9 | 49 | 19.1 | 15.1 | 55.9 | 43.6 | 14.4 | 12.3 |

$^a$ The percent of the relative amino acids.

TABLE 1. The physicochemical properties of the CNC caleosins.

| No. | Motif | Poaceae | Brassicaceae |
|---|---|---|---|
| 1 | SVLQQHVAF | [ST][VAP][LM]Q[QKG]H[VA]AF | SVLQQHVAF |
| 2 | D[RL]D[DG][ND]GI[IV]YPWETY | D[LR][DN][GNDK][DN]GI[IV]YP[WS]ET[YF] | D[LIM][DY][DG]NGIIYPWETY |
| 3 | NIH[KR][SG]KHGSDS | NIH[KR][GSDVA]KHGSDS[SGE]S | NIHKSKHGSDS[KRQ]S |
| 4 | GRFMPVN[FL]E[NL]IFSKY | GRF[MVD]P[VSE][NK]F[ED][NSEA]IF[SK]K[YH] | GRFMPVNLELIFSKY |
| 5 | LS[KR]EA[IV]RR[CM][FY]DGSLF | L[SAQH][KR][ED][AVIT][VIM]R[RGA][CMA][FY]DGSLF | LSKEAIRRCFDGSLF |

TABLE 2. The conserved motifs of the caleosins in Poaceae and Brassicaceae families. Motifs were predicted using the MEME Suite.

The plant families are marked with the symbols in the tree. Most of the sequences belonged to Poaceae and Brassicaceae families, and the rest were distributed among eight families: Asteraceae, Chlorellaceae, Euphorbiaceae, Fabaceae, Malvaceae, Rutaceae, Solanaceae, and Volvocaceae. As it is clear, each family occupies a separate category.

Since Poaceae indicated more diversity, the Greek numbers (I to IV) were applied to further investigations of the subclasses (Figure 1). While the I and II subclasses are thought to have diverged after the monocot-dicot split, the III and IV subclasses have emerged before the monocot-dicot split. A comparison of the CNC caleosins demonstrated that the C-terminal domain extended more in the IV, and then I, II and III subclasses, respectively. Another characteristic was higher pI value in the IV subclass. The net charge of the N-terminal domain probably has had changed over the evolution time because the III and IV subclasses and microalgae caleosins owned less negative charge on the N-terminal domains. In addition to the mentioned points, the percentage of aromatic amino acids (F+H+W+Y) was higher in I and II subclasses.

3.3. **Conserved motif in the calseosins.** Motif study revealed five conserved regions in the caleosins (Table 2). According to the motif library of GenomeNet, all of the recognized motifs belonged to the caleosin family (pfam05042). Interestingly, a motif investigation in CDART demonstrated that these motifs can also be found in phosphatase of *Aspergillus* sp. (GI: 846917092 and 849274718), ubiquitin-associated protein 2 of *Heterocephalus glaber* (GI: 351713888), calcineurin (GI: 695044459 and 927143482), and some unknown sequences with formate-dependent nitrite reductase activity (GI: 4148 65577, 590659846, 823247273 and 823247271).

Since most of the caleosins belonged to Poaceae and Brassicaceae families (Figure 1), the sequence motifs of these two families were presented in Table 2. As it is clear, members of Brassicaceae were highly conserved, while a great diversity was observed in Poaceae. It could represent the long-term evolution in the Poaceae family.

Each motif was usually found in a particular domain (Table 2) [9]. Compared to motif 5 which mostly stands as a caleosin sign on the C-terminus, motif 1 and 2 along with a conserved F residue between them are placed predominantly on the N-terminal domain. However, A0A061EME7 from *Theobroma cacao* includes one additional repetition of motif 1 and 2 on the central domain. Motif 1 is a heme-binding fragment with a highly conserved histidine in SVLQQHVAF sequence that is essential for peroxygenase activity [21]. Motif 2 or EF-hand was first discovered in the calcium-binding proteins of bacteria, and then in the fungal caleosins and lectins of *Psathyrella velutina* [11]. This motif holds 26 amino acids with the consensus sequence of Dx[DN]xDG that makes a helix-loop-helix structure. Aspartate or glutamate residues of the motif provide
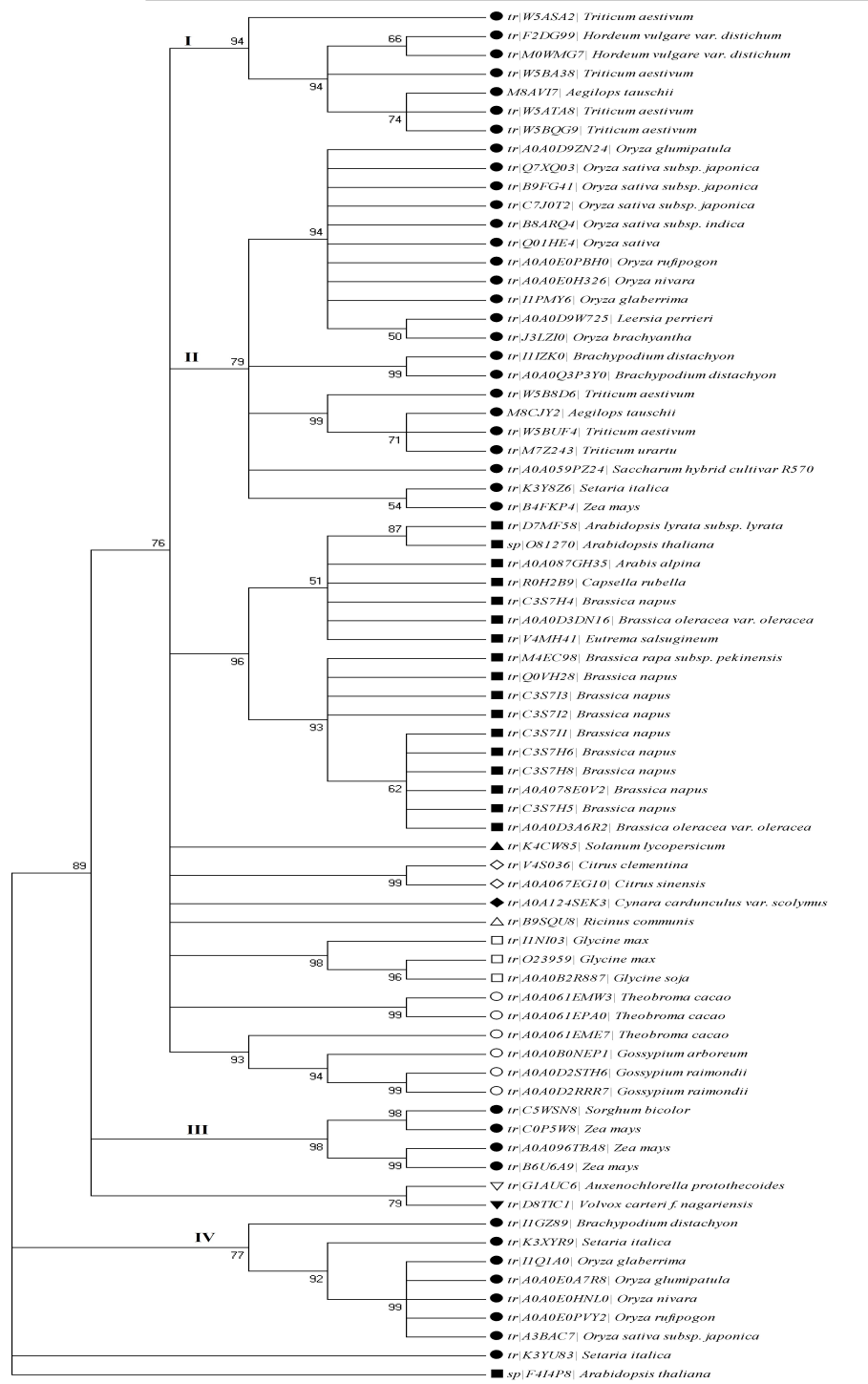
FIGURE 1. Neighbor-joining phylogenetic tree of the caleosins. The numbers at the nodes indicate bootstrap values based on 1000 replications. ●, ■, ▲, ▼, ♦, ○, □, △, ▽, ◇ represent caleosins from Poaceae, Brassicaceae, Solanaceae, Volvocaceae, Asteraceae, Malvaceae, Fabaceae, Euphorbiaceae, Chlorellaceae, and Rutaceae families, respectively.

oxygen ligands for the calcium-binding site [17; 20]. Although the exact role of this motif has remained unknown, it could be involved in the early calcium-induced signaling pathway during the plant defense responses [23].

Generally, motif 3 and 4 are following the proline knot motif on the central domain. However, they were also detected in the transmembrane region and the C-terminus of A0A061EME. Moreover, A0A067EG10 and V4S036 from citrus consists of a less conserved repetition of motif 3 and 4 on the N-terminal domain. Although the proline knot is the best-known motif in the oil body-membrane proteins, it was not verified by MEME probably because of high substitutions.

3.4. **A model sequence for caleosins.** A consensus sequence was extracted to describe the protein properties such as hydrophobicity, secondary structure, post-translational modifications, and subcellular locations. For this purpose, sequences with low identities, such as F4I4P8 were excluded from the CNC group. Then, a sequence model was developed by WebLogo (Figure 2A). Most parts of the sequences were conserved; however, the N-terminal domain of H-caleosins [22] showed diversity. As it is clear in Figure 2A, acidic and basic amino acids have positioned intermittently. Consequently, an efficient system is constructed in which the negative charges prevent from the oil body aggregation by the electrostatic repulsion, while the positive charges bind to the negatively charged phospholipid surface and clamp the terminal domains on the oil body surface. This feature, along with the amphipathic nature of caleosins (Figure 2B), avoids detaching of the terminal domains from the oil body surface [26].

A domain analysis verified the three structural domains with a probability of 90% in the sequence model (Figure 2C). However, no signal peptide was detected in the CNC caleosins. The five motifs were specified in the sequence (Figure 3). As can be seen in Figure 3, the alpha helix with 35.5%, after the random coil (56.03%), shapes most of the caleosin sequences. Like oleosins [27], the presence of the alpha-helix structure was predominant in the C-terminal of caleosins.

*In silico* study identified the phosphorylation and ubiquitination sites (Figure 3). Although the possibility of acetylation was estimated less than 0.5 by NetAct, a mass spectrometry analysis by Lin et al. (2005) has proved the N-terminal acetylation of a caleosin. This modification leads to protein stability by inhibiting the ubiquitin-proteasome degradation system [15; 18]. Notably, the examined caleosin is not among the CNC group due to the lack of the N-terminal domain. In addition to the mentioned above, N- and O- glycosylation was also predicted for some caleosins [18].
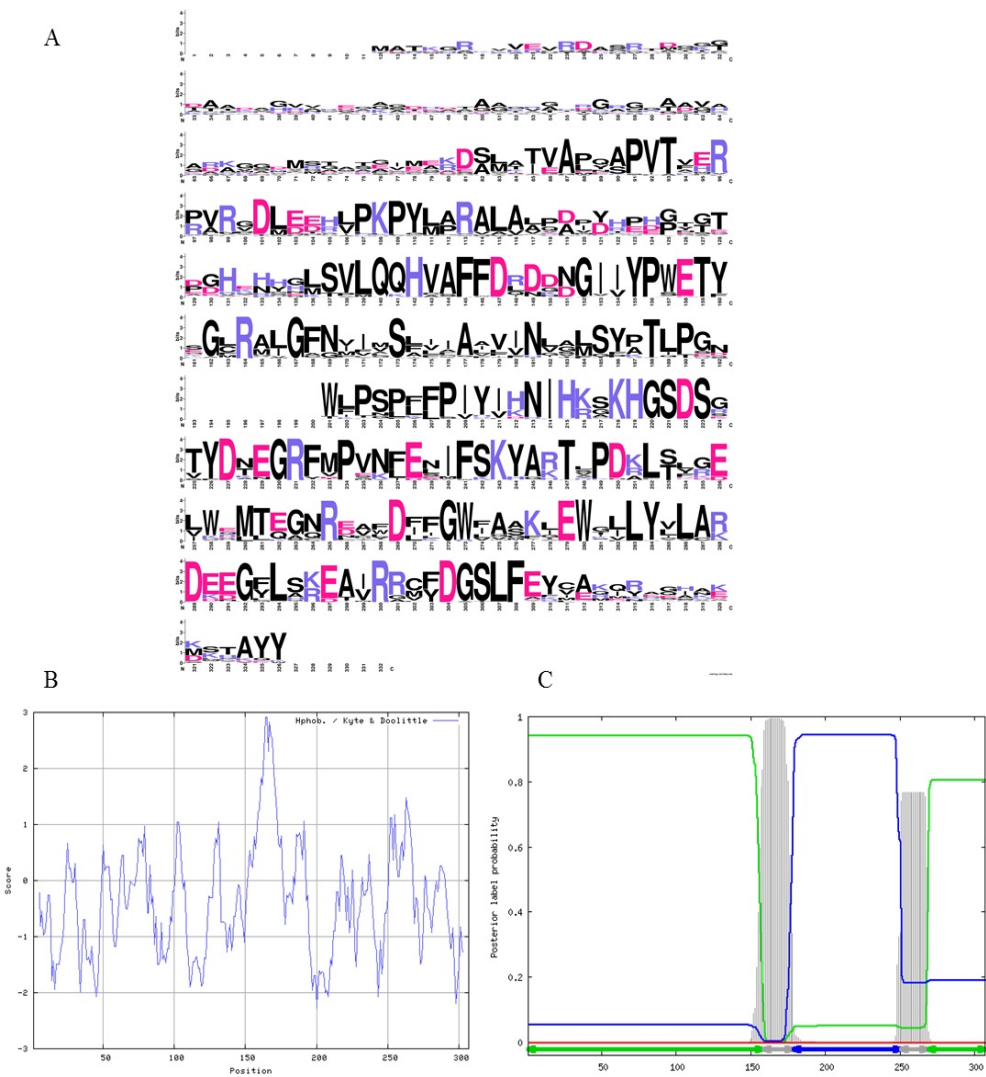
F. Saadat, H. Alizadeh, S. H. Razavi

A



B

C

FIGURE 2. A model for the CNC caleosins. A) Sequence logo of the CNC caleosins. Blue and pink words represent the basic and acidic amino acids, respectively. B) Hydrophobicity plot of the consensus sequence. C) Transmembrane topology of the consensus sequence in which grey parts represent the transmembrane regions, and blue, green and red lines show non-cytoplasmic, cytoplasmic and signal peptide regions, respectively.

**N-terminal Domain:**

MA**T**KGRRVVEVRDA**S**RTRGGGDAADAGVVRE
ccccccceeeeeecccccccccccchhhhhhh
GAGDHDTAAGRGHRGRGGAAVAAR**K**GGDM**S**T
cccccccccccccccchhhhhhhccccccee
A**T**GIMEKD**S**LA**T**VAPQAPVTVERPVRGDLEE
eehhhhhhhhhcccccccceeccccccccccc
HLPKPYLARALALPDP**Y**HPHG**T**G**T**PGHEHHG
cccchhhhhhhcccccccccccccccccccc
L**S**VLQQHVAFF DRDDNGIIYPWE**T**Y SGLRAL
hhhhhhhhhhhcccccceeecccccccccccc
GFN
cch

**Central Domain:**

PGNWLP**S**PFFPI**Y**IH NIHKSKHG**S**D**S**G
cccccccccccceeeecccccccccccccc
**TY**DNE GRFMPVNFENIF**S**K**Y**AR**T**
cccccccccccccccccchhhhhcc
LPDKL**S**LGELWEMTEGNREAFD
ccccccccchhhhhhhcccchhhh

**C-terminal Domain:**

RDEEGF **L**S**KEAIRRCFDG**S**LF**EYC
cchhhhhhhhhhhhhcccccchhhh
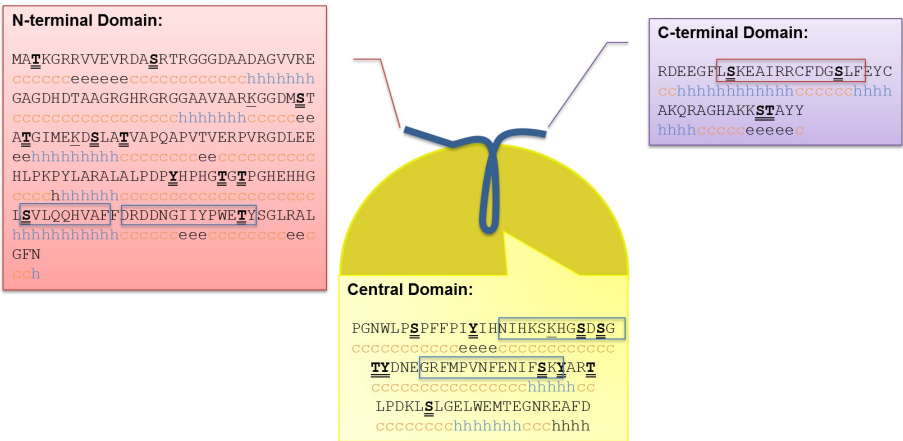AKQRAGHAKK**ST**AYY
hhhhcccccceeeeec

FIGURE 3.  The consensus sequence of the CNC caleosins.  The
sequence of each structural domain mentioned separately.  In
addition, the secondary structure was shown in the sequence
with c, h, and e represent coil, helix, and beta-sheet structure,
respectively.  Each motif has been placed inside a box.  Ubiqui-
tination and phosphorylation sites have been displayed by one
and two underlines, respectively.

## 4. CONCLUSION

In the current study, a comprehensive investigation was carried out on the
reported plant caleosins which possess the CNC structure. In agreement with
other papers, about 30% of a caleosin sequence has devoted to the central
domain, which was nearly half-length of the terminal domains [17]. The long
terminal domains help caleosins build smaller droplets.

Comparing the structural features of the caleosins in two major plant fami-
lies, Brassicaceae and Poaceae, made a proper representative for dicots and
monocots plants, respectively.  Accordingly, members of Brassicaceae were
highly conserved not only in the motif sequences but also in the average length
and charge of domains (data not shown). In contrast, there was a high diversity
among the members of Poaceae that could be a sign of their long evolutionary
history.  This variation created several clades for Poaceae in the phylogenetic
tree.  Further analysis also showed a shorter N-termini, more positive charge,
higher pI value, and less aromatic residues in the ancient species. Based on
the results, the negative charge of the terminal domains and the length of the
hydrophobic region have changed over evolutionary time. However, its confir-
mation requires more investigation.

## References

1. T. L. Bailey, N. Williams, C. Misleh, W. W. Li, MEME: Discovering and Analyzing DNA and Protein Sequence Motifs, *Nucleic Acids Research*, **34**, (2006), W369-373.

2. S. C. Bhatla, V. Kaushik, M. K. Yadav, Use of Oil Bodies and Oleosins in Recombinant Protein Production and other Biotechnological Applications, *Biotechnology Advances*, **28**(3), (2010), 293-300.

3. V. Brendel, P. Bucher, I. R. Nourbakhsh, B. E. Blaisdell, S. Karlin, Methods and Algorithms for Statistical Analysis of Protein Sequences, *Proceedings of the National Academy of Sciences of the United States of America*, **89**(6), (1992), 2002-2006.

4. M. C. M. Chen, C. L. Chyan, T. T. T. Lee, S. H. Huang, J. T. C. Tzen, Constitution of Stable Artificial Oil Bodies with Triacylglycerol, Phospholipid, and Caleosin, *Journal of Agricultural and Food Chemistry*, **52**(12), (2004), 3982-3987.

5. C. J. Chiang, L. J. Lin, C. J. Chen, Caleosin-based Nanoscale Oil Bodies for Targeted Delivery of Hydrophobic Anticancer Drugs, *Journal of Nanoparticle Research*, **13**(12), (2011), 7127-7137.

6. C. J. Chiang, S. C. Lin, L. J. Lin, C. J. Chen, Y. P. Chao, Caleosin-assembled Oil Bodies as a Potential Delivery Nanocarrier, *Applied Microbiology and Biotechnology*, **93**(5), (2012), 1905-1915.

7. C. Combet, C. Blanchet, C. Geourjon, G. Delage, NPS@: Network Protein Sequence Analysis, *Trends in Biochemical Sciences*, **25**, (2000), 147-150.

8. G. E. Crooks, G. Hon, J. M. Chandonia, S. E. Brenner, WebLogo: A Sequence Logo Generator, *Genome Research*, **14**(6), (2004), 1188-1190.

9. Y. Fan, A. Ortiz-Urquiza, T. Garrett, Y. Pei, N. O. Keyhani, Involvement of a Caleosin in Lipid Storage, Spore Dispersal, and Virulence in the Entomopathogenic Filamentous Fungus, Beauveria Bassiana, *Environmental Microbiology*, **17**(11), (2015), 4600-4614.

10. L. Y. Geer, M. Domrachev, D. J. Lipman, S. H. Bryant, CDART: Protein Homology by Domain Architecture, *Genome Research*, **12**(10), (2002), 1619-623.

11. A. P. Gonalves, J. Monteiro, C. Lucchi, D. J. Kowbel, J. M. Cordeiro, P. Correia-de-S, D. J. Rigden, N. L. Glass, A. Videira, Extracellular Calcium Triggers Unique Transcriptional Programs and Modulates Staurosporine-induced Cell Death in Neurospora Crassa, *Microbial Cell*, **1**, (2014), 289-302.

12. P. -L. Jiang, J. T. C. Tzen, Caleosin Serves as the Major Structural Protein as Efficient as Oleosin on the Surface of Seed Oil Bodies, *Plant signaling & behavior*, **5**(4), (2010), 447-449.

13. L. Kll, A. Krogh, E. L. L. Sonnhammer, Advantages of Combined Transmembrane Topology and Signal Peptide Prediction–the Phobius Web Server, *Nucleic Acids Research*, **35**, (2007), W429-432.

14. J. Kyte, R. F. Doolittle, A Simple Method for Displaying the Hydropathic Character of a Protein, *Journal of Molecular Biology*, **157**(1), (1982), 105-132.

15. L. J. Lin, P. C. Liao, H. H. Yang, J. T. C. Tzen, Determination and Analyses of the N-termini of Oil-body Proteins, Steroleosin, Caleosin and Oleosin, *Plant Physiology and Biochemistry*, **43**(8), (2005), 770-776.

16. D. Murphy, The Biogenesis and Functions of Lipid Bodies in Animals, Plants and Microorganisms, *Progress in Lipid Research*, **40**, (2001), 325-438.

17. Z. Purkrtov, T. Chardot, M. Froissard, N-terminus of Seed Caleosins is Essential for Lipid Droplet Sorting but not for Lipid Accumulation, *Archives of Biochemistry and Biophysics*, **579**, (2015), 47-54.

18. Z. Purkrtova, S. dAndrea, P. Jolivet, P. Lipovova, B. Kralova, M. Kodicek, T. Chardot, Structural Properties of Caleosin: A MS and CD Study, *Archives of Biochemistry and Biophysics*, **464**(2), (2007), 335-343.

19. P. Radivojac, V. Vacic, C. Haynes, R. R. Cocklin, A. Mohan, J. W. Heyen, M. G. Goebl, L. M. Iakoucheva, Identification, Analysis, and Prediction of Protein Ubiquitination Sites, *Proteins*, **78**(2), (2010), 365-380.

20. V. S. Reddy, A. S. N. Reddy, Proteomics of Calcium-signaling Components in Plants, *Phytochemistry*, **65**(12), (2004), 1745-1776.

21. Y. Shen, M. Liu, L. Wang, Z. Li, D. C. Taylor, Z. Li, M. Zhang, Identification, Duplication, Evolution and Expression Analyses of Caleosins in Brassica Plants and Arabidopsis Subspecies, *Molecular genetics and genomics*, **291**(2), (2016), 971-988.

22. Y. Shen, J. Xie, R. -D. Liu, X. -F. Ni, X.-H. Wang, Z. -X. Li, M. Zhang, Genomic Analysis and Expression Investigation of Caleosin Gene Family in Arabidopsis, *Biochemical and biophysical research communications*, **448**(May), (2014), 365-371.

23. T. L. Shimada, I. Hara-Nishimura, Leaf Oil Bodies are Subcellular Factories Producing Antifungal Oxylipins, *Current Opinion in Plant Biology*, **25**, (2015), 145-150.

24. W. Song, Y. Qin, Y. Zhu, G. Yin, N. Wu, Y. Li, Y. Hu, Delineation of Plant Caleosin Residues Critical for Functional Divergence, Positive Selection and Coevolution, *BMC evolutionary biology*, **14**(1), (2014), 124.

25. K. Tamura, D. Peterson, N. Peterson, G. Stecher, M. Nei, S. Kumar, MEGA5: Molecular Evolutionary Genetics Analysis Using Maximum Likelihood, Evolutionary Distance, and Maximum Parsimony Methods, *Molecular Biology and Evolution*, **28**(10), (2011), 2731-2739.

26. T. Wahlroos, J. Soukka, A. Denesyuk, P. Susi, Amino-terminus of Oleosin Protein Defines the Size of Oil Bodies -Topological Model of Oleosin-oil Body Complex, *Journal of Plant Biochemistry and Physiology*, **3**, (2015), 1-6.

27. D. Zweytick, K. Athenstaedt, G. Daum, Intracellular Lipid Particles of Eukaryotic Cells, *Biochimica et Biophysica Acta*, **1469**, (2000), 101-120.